(19) World Intellectual Property Organization International Bureau





(43) International Publication Date 27 December 2001 (27.12.2001)

PCT

C12N

(10) International Publication Number WO 01/98454 A2

(51) International Patent Classification7: (21) International Application Number: PCT/IB01/02050 (22) International Filing Date: 25 April 2001 (25.04.2001) (25) Filing Language: English (26) Publication Language: English

(30) Priority Data: 25 April 2000 (25.04.2000) 60/199,380

(63) Related by continuation (CON) or continuation-in-part (CIP) to earlier application: US 60/199.380 (CIP) Filed on 25 April 2000 (25.04.2000)

(71) Applicant (for all designated States except US): GER-MAN HUMAN GENOME PROJECT [DE/DE]; Fraunhofer Patentstelle, Leonrodstrasse 68, 80636 Munich (DE).

(72) Inventor; and

(75) Inventor/Applicant (for US only): WIEMANN, Stefan

[DE/DE]; Grosse Lachstrasse 30a, 69207 Sandhausen

(81) Designated States (national): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CR, CU, CZ, DE, DK, DM, DZ, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, TZ, UA, UG, US, UZ, VN, YU, ZA, ZW.

(84) Designated States (regional): ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).

Published:

without international search report and to be republished upon receipt of that report

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

(54) Title: HUMAN DNA SEQUENCES

(57) Abstract: Novel human cDNA sequence of a clones, the encoded protein sequence of a clones, antibodies and variants thereof, are provided. The disclosed sequence of a clones find application in a number of ways, including use in profiling assays. In this regard, various assemblages of nucleic acids or proteins are provided that are useful in providing large arrays of human material for implementing large-scale screening strategies. The disclosed sequence of a clones may also be used in formulating medicaments, treating various disorders and in certain diagnostic applications.

> Applicants: Paz Einat et al. Serial No.: 10/618,408 Filed: July 11, 2003

Exhibit 5

WO 01/98454

10

15

20

25

30

35



HUMAN DNA SEQUENCES

Background of the Invention

Current methods for testing pharmacological substances rely on a three-stage testing approach to drug development. First, candidate compounds are typically screened in some sort of in vitro system, like inhibition of cancer cell growth. Candidates are then tested in an animal model, as a first approximation of systemic effects, including efficacy and toxicity. Compounds that still show promise after these initial in vivo screens, finally are tested in humans. Again, human testing typically occurs in three phases: toxicity; preliminary efficacy; and efficacy. The entire process can take more than a decade and cost hundreds of millions of dollars. Aside from the monetary costs and protracted time scale, moreover, current testing regimes waste the lives of countless laboratory animals and needlessly endanger the lives of human subjects.

A need exists, therefore, for more sophisticated drug screening techniques that can be done rapidly in vitro. These screening techniques ideally will be reflective of systemic and/or organ-specific responses, so that they provide a reliable indicator of action in a human body. Current techniques, however, tend to utilize only a single or limited number of markers, thus answering only very simple questions that are of questionable medical import. For example, a typical in vitro assay may ask whether a lead compound binds a particular receptor, which has been implicated in a certain disorder. It is presumed that such binding is indicative of therapeutic usefulness, but it does not even purport to address systemic effects.

Not only are screening techniques for efficacy inadequate: the available toxicity screens likewise are inadequate. Toxicity: on a first level: is usually measured by animal testing. Aside from the complications related to in vivo versus in vitro testing: such screens are insufficient because of differences in metabolism: uptake: etc.: relative to humans.

Thus, improved methods would be not only be in vitro-based, they would also be more "human."

With the increasing miniaturization of screening assays and the growing availability of targets for pharmaceutical 5 intervention, there is increasing interest in developing arrays containing large numbers of these targets that can be assayed simultaneously. If such an array contains a large enough population of targets, it can be used to essentially mimic the systemic response. In other words, the array becomes an in vitro 10 surrogate for the human body. The more refined the array, the more accurate the predictive capability. In theory, an array could be constructed that can detect all of the known human expression products simultaneously, thereby, providing a very reliable indicator of the human response to a given compound. These arrays offer advantages over the present in vitro screening 15 systems in that they can assay large numbers of responses simultaneously. They are superior to animal testing because they are more "human" and, thus, more predictive of human responses.

In order to construct such arrays, however, the field is in need of further human targets. Advantageously, such targets will be provided with additional physiologically relevant information, such as whether the target is expressed in a particular tissue and whether it is related to a known functional class of targets. In this way, the artisan can focus as needed, for example, on 25 tissue-specific effects or target class-specific effects, thereby providing information useful in evaluating efficacy and/or toxicity.

20

30

35

In addition to a need for pharmacological screening targets, there is a need for further pharmacological substances. These substances can be used in the formulation of medicinal compositions and in treating a wide variety of disorders.

The present invention responds to the aforementioned and other needs in the field by providing a population of novel targets useful, inter alia, in the profiling and medicinal contexts described above.

Summary of the Invention

It is an object of the invention, therefore, to provide a set of human cDNA clones. Further to this object, the invention provides sequences of human cDNA clones that were isolated from libraries generated from different human tissues.

It is another object of the invention to provide assemblages of targets useful in profiling matrices for screening pharmacological test compounds. According to this object, assemblages comprising different populations of human nucleic acids, proteins and antibodies are provided. In different embodiments, cDNA library-specific assemblages and target-family-specific targets are provided.

10

15

20

25

30

35

It is a further object of the invention to provide a database of human nucleotide and protein sequences. Further to this object, novel human nucleotide and protein sequences are provided in electronic form. In one embodiment, one or more of these sequences is provided in a searchable database.

It is still another object of the invention to provide biologically active target molecules useful in treating or detecting human disorders. Further to this object, the invention provides nucleic acid and protein molecules that have the capacity to affect disease etiology or symptoms or correlate with known disease states. Also further to this object, a database is provided which comprises the disclosed molecules in electronic form.

Detailed Description

The invention results from a need in the art for new human nucleic acids and proteins. This need arises in several contexts. First, there is a need to identify targets for therapeutic intervention. Second, there is a need to identify molecules that may be adversely affected in a therapeutic context, thereby resulting in toxicity. Knowledge of these molecules will aid in the design of new medicaments with enhanced efficacy and decreased toxicity. Finally, the need encompasses human nucleic acids and proteins that have medicinal applicability in their own right.

In view of these needs, the present inventors set out to isolate and sequence human cDNAs from tissue-specific libraries.

-3-

In this way, they represent subsets of molecules likely to be targets for therapeutic intervention or for avoiding toxicity. In addition, the inventors divided the molecules into various subcategories, based on suspected functionality, structural similarity etc, which are of interest from a pharmacological perspective.

GENERAL DESCRIPTION OF THE INVENTIVE MOLECULES

10

15

20

25

30

35

The present invention provides novel polynucleotide molecules that in some instances have similarities with known molecules. The inventive DNAs were cloned from five different human cDNA libraries. In addition to these DNA molecules the invention provides their protein translations and antibodies derived from them. The inventive DNA and protein sequences are show individually in the Description of the Sequences. The inventive nucleic acids also include the complements of the DNA sequences provided in the Description of the Sequences as well as their RNA counterparts. Methods of producing the molecules also are provided. Further, the invention provides methods for detecting all or part of the molecules and of detecting polynucleotides encoding all or part of the molecules.

The inventive molecules derive from five cDNA libraries: human fetal brain; human fetal kidney; human melanoma; human testis; and human amygdala. For convenience, each sequence bears a designation that indicates from which library it is derived. In particular, these designations are: "hfpbr" for human fetal brain; "hfkd" for human fetal kidney; "hmel" for human melanoma; "htes" for human testis; and "hamy" for human amygdala. The individual libraries were constructed and screened as described below in the examples.

The protein and DNA molecules of the invention are variously described herein as "target" molecules or "inventive" molecules. The sequences and other information pertinent to the nucleic acid and protein molecules of the invention are shown below in the Description of the Sequences.

Description of the Sequences
Key to the Description of the Sequences

The desctiptions below provide the coding sequences of the inventive cDNAs, as well as the protein sequences and other useful information, as set out herein.

Grouping

The clones were assigned to the following sixteen functional and/or tissue-derived groups:

- 1. Amygdala derived
- 10
- 2. Cell Cycle 3. Cell Structure and Motility
 - 4. Differentiation/Development
 - 5. Intracellular Transport and Trafficking
 - **b** Melanoma derived
- 15 7. Metabolism
 - 8. Nucleic Acid Management
 - 9. Signal Transduction
 - 3D-Transmembrane Protein
 - 11. Transcription Factors
- 20 15. Brain derived
 - 13. Kidney derived
 - 14. Mammary Carcinoma derived
 - 15-Testes derived
 - **16** -Uterus derived

25

Description of Clone Files

The individual clone files are structured in the same pattern. The Sections are separated by paragraphs.

30 1. Clone Name

> The clone names are deciphered with reference to the following example:

> > DKFZphfkd2_3kl, wherein the code represents:

- producer of library ("DKFZ") (for convenience, this reference may be eliminated)
 - a "p" for "plasmid cDNA library" (for convenience, this reference may be eliminated)
 - library name (e.g. hfbr = human fetal brain; hfkd = human fetal kidney; hmel = human melanoma; htes = human testis: hamy = human amygdala)
 - an underscore ("_") to separate library information from plate information
 - plate number (e.g. "3")
 - plate coordinates (letter first; e.g. "kl2")

45

40

35

2. Group

3. Introduction

short review of the similarities, function of the protein and possible applications

5 4. Short Information

specifications about the cDNA (who sequenced, completeness of the cDNA, similarity, who sequenced, chromosomal localisation, length of cDNA, localisation of poly A tail and polyadenylation signal)

- 10 5. cDNA-Sequence
 - 6. BLASTn Results

search results of blasting the cDNA sequence against all public databases

15

7. Medline Entries

information about genes/proteins similar to the novel cDNA (if available)

- 20 8. Putative Encoded Protein Information specifications about the encoded protein (ORF: length and localisation of the reading frame)
 - 9. Protein Sequence

25

10. BLASTp Results

search results of blasting the protein sequence against all public databases

30 11. Pedant Information

output of fully automated annotation: summarises peptide information, homologies, patterns as follows:

[Length]

- length of the protein = number of amino acid residues
 - molecular weight of the protein $\ensuremath{\mathtt{TpII}}$

- isoelectric point

entry point to the database.

ELOMOLI

shows protein with closest similarity to the cDNA-encoded protein

EFUNCATI

 functional information according to a catalogue developed by Munich Information center for Protein Sequences (MIPS)

EBLOCK23

10

15

20

5 ·

- Blocks are multiply aligned ungapped segments corresponding to the most highly conserved regions of proteins. The blocks for the Blocks Database are made automatically by looking for the most highly conserved regions in groups of proteins documented in the Prosite Database. The Prosite pattern for a protein group is not used in any way to make the Blocks Database and the pattern may or may not be contained in one of the blocks representing a group. These blocks are then calibrated against the SWISS-PROT database to obtain a measure of the chance distribution of matches. It is these calibrated blocks that make up the Blocks Database. The WWW versions of the Prosite and SWISS-PROT Databases that are used on this server are located at the ExPASy World Wide Web (WWW) Molecular Biology Server of the Geneva University Hospital and the University of Geneva. World Wide Web URL http://blocks.fhcrc.org/blocks/about_blocks.html/ is the
- 25
- here Blocks segments found in the analysed protein sequences are displayed $\ensuremath{\mathtt{ESCOPI}}$

30

35

Nearly all proteins have structural similarities with other proteins and, in some of these cases, share a common evolutionary origin. The scop database provides a detailed and comprehensive description of the structural and evolutionary relationships between all proteins whose structure is known, including all entries in Brookhaven National Laboratory's Protein Data Bank (PDB). It is available as a set of tightly linked hypertext documents which make the large database comprehensible and accessible.

5.

10

15

20

25

30

35

In addition, the hypertext pages offer a panoply of representations of proteins, including links to PDB entries, sequences, references, images and interactive display systems. World Wide Web URL http://scop.mrclmb.cam.ac.uk/scop/ is the entry point to the database. Existing automatic sequence and structure comparison tools cannot identify all structural and evolutionary relationships between proteins. The scop classification of proteins has been constructed manually by visual inspection and comparison of structures, but with the assistance of tools to make the task manageable and help provide generality. Proteins are classified to reflect both structural and evolutionary relatedness. Many levels exist in the hierarchy, but the principal levels are family, superfamily and fold. The exact position of boundaries between these levels are to some degree subjective. Scop evolutionary classification is generally conservative: where any doubt about relatedness exists, we made new divisions at the family and superfamily levels.

 - here SCOPE segments found in the analysed protein sequences are displayed
 LECI

ENZYME is a repository of information relative to the nomenclature of enzymes. It is primarily based on the recommendations of the Nomenclature Committee of the International Union of Biochemistry and Molecular Biology (IUBMB) and it describes each type of characterized enzyme for which an EC (Enzyme Commission) number has been provided. World Wide Web URL http://www.expasy.ch/enzyme/ is the entry point to the database.

- here EC-number and name of enzymes with similarity to the analysed protein sequences are displayed $\ensuremath{\mathtt{EPIRKWI}}$
- functional information according to the Protein Information Resource (PIR) database catalogue developed by Munich Information Center for Protein Sequences (MIPS), the National Biomedical Research Foundation (NBRF) and the International Protein Information Database in Japan (JIPID). ESUPFAMD

- information according to the Protein Information Resource (PIR) database catalogue of protein superfamilies developed by Munich Information Center for Protein Sequences (MIPS), the National Biomedical Research Foundation (NBRF) and the International Protein Information Database in Japan (JIPID).

IPROSITE

5

20

25

30

35

please refer to 12. PROSITE Motifs

[PFAM]

10 please refer to 13. PFAM Motifs

- overall 2dimensional folding information
- 3D indicates that the proteins is similar to a protein of which a 3 dimensional structure is known
- overall structural information

The last PEDANT-block depicts information about the folding structure of the protein generated by PREDATOR.

PREDATOR is a secondary structure prediction program. It takes as input a single protein sequence to be predicted and can optimally use a set of unaligned sequences as additional information to predict the query sequence. The mean prediction accuracy of PREDATOR is LAX for a single sequence and 75% for a set of related sequences. PREDATOR does not use multiple sequence alignment. Instead, it relies on careful pairwise local alignments of the sequences in the set with the query sequence to be predicted.

World Wide Web URL http://www.emblheidelberg.de/argos/predator/predator_info.html is the entry point to the database.

- H = helix, E = extended or sheet, _ = coil, T = transmembrane, B = beta
- x indicates a low-complexity region with repeat-like structure which is omitted in all BLAST searches

12. PROSITE Motifs

PROSITE is a database of protein families and domains. It consists of biologically significant sites, patterns and profiles that help to reliably identify to which known protein family (if

any) a new sequence belongs. World Wide Web URL http://www.expasy.ch/prosite/ is the entry point to the database. A description of the prosite consensus patterns is provided herein, after the description of the individual sequences.

5

10

13. PFAM Motifs

PFAM (protein families) is a large collection of multiple sequence alignments and hidden Markov models covering many common protein domains. World Wide Web URL http://www.sanger.ac.uk/Pfam/is the entry point to the database.

In the charts below, the groups of sequences are listed, and the description of the individual clones follows.

Group Amygdala derived

÷.	1	1	 -	T	т	г	
Group	amygdala	amygdala	amygdala derived	amygdala derived	Amygdala derived	amygdala derived	amygdala derived
Function	No informative BLAST results: No predictive prosite, pfam or SCOP motife.	No informative BLAST results; No predictive prosite, pfam or SCOP motife	The novel protein contains a PROSITE ASP_PROTEASE motif and seem to be expressed Ubiquitously. No informative BLAST results No predictive procite, of an or stop motifs.	A similar CDNA encoding a protein of the same length was identified in sheep. This amygdala protein shows a strong signal sequence, which indicates that it is a secreted protein. The new protein belongs to a protein family, which was designated carbonic anhydrass-related protein XI (CA-RP XI), encoded by (All (human) and Call (mouse, rel). Despite potentially inactivating changes in the active-site residues, CA-RP XI is evolving very slowly in mammals, a property indicative of an important function, which has also been observed in the two other "acatalytic" (A isoforms, CA-RP VIII and CA-RP X. No informative BLAST results: No predictive prosite, pfam or SCOP motife	Pecanex is a maternal-effect neurogenic gene, involved in differentiation processes in the developing central nervous system. DKFZphamy2_24kl5 seems to be expressed ubjquitiously.	No informative BLAST results: No predictive prosite, pfam or SCOP motife	Most ESTs are derived from brain and pancreas No informative BLAST results: No predictive prosite, pfam or SCOP motife.
Homology	Without similarity to known proteins	weak similarity to F41E6.3 of Caenorhabditis elegans	amy2_13g14 without similarity to known proteins	amy2_lbel4 similar to carbonic anhydrasa- related proteins	amy2_24k15 weak similarity to pecanex of Drosophila melanogaster.	without similarity to known proteins	amy2_zil7 without similarity to known proteins
E CloneID DEFEDS	amy2_12g7	amy2_leil	amy2_13g19	amy2b60l4	amy2_24k1s	amy2_2əl3	amy2_2i17

Group Brain derived

DKFZph	Homology				Punction			dnozb
fbr2_78d18	br2_75d16 weak similarity to a human	-	No informative BLAST results: No predictive prosite, plam or SCOP motife.	Isults: No	predictive prosite.	Dfam or S	COP motife.	brain
	putative mitogen-activated protein	otein						derived
	kinase kinase kinase							
fbr2_78e38	bre_78el8 without similarity to known		The mRNA is differentially polyadenylated.	illy polyac	senvlated.			brain
	proteins.		No informative BLAST ru	sults: No	predictive prosite.	ofam or S	CoP motife.	derived

Group cell cycle

		
A Quotable of the second of th	PARE-TE is a p53 responsive gene. The protein is predominantly expressed in brain, cell cycle breast and kidney and may represent a potential novel regulator of cellular growth. Isoforms are differentially induced by genotoxic stress (UV, gamma-Irradiation and cutokoxic druns) in a m51-denomber manner.	15 16
lg.	2	20
	8	Ce
	يَ	The stromal interaction molecular 1 gene (SIM1) encodes a type I trans-membrane cell cyle protein of unknown function, which induces growth arrest and degeneration of the human tumor cell lines 6401 and RD but not HBLMO and Calu-6, suggesting a role in the pathogenesis of rhabdomyosarcomas and rhabdoid tumors. There is also strong similarity to a Mus musculus stromal cell protein, which selectively increases interleukin 7-dependent proliferation of pre-8 cells. The novel protein contains 1 transmembrane domain.
	4	roll roll ron ses
	4 1 2 E	on on o
	asse cell	rans rati stin als als
	z P	I til
1.0	1y tor	ype deg deg here ecti
	gula str	Sel sel
12	domi 1 re 0xic	rest a Ca a Ca a Ca Th
	nove not	1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2
	al y	B to be constant
110	tein enti	S STILL S S STILL S S STILL S STILL S S STILL S S STILL S S STILL S S
E .	Page 1	t duce duce and call
;÷	The I	l d h in h in h in hal
18	PA26-TE is a p53 responsive gene. The protein is predominantly expressed in breast and kidney and may represent a potential novel regulator of cellular growth. Isoforms are differentially induced by genotic stress (UV. gamma-freediation and cutotoxic frunts) in a p52-denominal manner.	The strongl interaction molecular in gene (STIA) encodes a type I trans-membrane protein of unknown function, which induces growth arrest and degeneration of the human tumor cell lines 6401 and RD but not HBL100 and Calu-6, suggesting a role is the pathogenesis of rhabdomyosarcomas and rhabdoid tumors. There is also strong similarity to a Mus musculus stromal call protein, which selectively increases interleukin 7-dependent proliferation of pre-B calls. The novel protein contains transmembrane domain.
1. M	repr	Diecu Di ar Myos us s
27.	asiv iff	n mo ctio 640 abdo scul
6.5	and and	Ctio fun ines r rh s mu
,	Sa r	tera nown 11 1 12 0 15 0 depe
	P S S S S S S S S S S S S S S S S S S S	unk unk ce enes to to
	and Is	no de
	E-T	The stromal interacti protein of unknown fu human tumor cell line the pathogenesis of r similarity to a Mus m linterleukin 7-depain.
	PAG Bre	trattur tra
1. /		
	6-T2	Ē
6	PAE	ELS
Homology	uman	uman
留	0. 1	음
-	ty 1	<u>t</u>
	lari ein.	lari
"	Similarit protein.	Simi
	my2_121m2 Similarity to human PA28	amy2_2464 Similarity to human STI
CloneID DKFZph	ัก สา	246
DXFZ	my2_	my 2,
لتا	ıŭ .	

Group Cell structure and motility

	The second secon		
CloneID	Homology	Constitution of the second sec	dronb
ZJfJ9 hi bi mR	myz_lzlfly high similarity to a Rat ank binding glycoprotein-l rel mRNA.	Rat ankyrin Ankyrin binding glycoproteins play a role in neural cell adhesion and in prosate cell n-l related twor cell transformation. DKFZphamy2_b2lf19 is expressed in brain, uterus and struct prostate above average	cell structure
ibbs si	es3_lbb5 similarity to various tropomyosins.	Tropomyosins play regulatory roles in cellular structure and transport.	cell structure

Group Differentiation/Development

CloneID DKFZph	Homology	The second of th	dnoap
amy2_1i24	partial similarity to rattus norvegicus Notch2 protein	euronal ot only by the postnatal	differentiat ion/developm ent
amyz_lji9	amy2_1j19 high similarity to the allograft inflammatory factor-1 of Cyprinus carpio.	Allograft inflammatory factor-1 (AIF-19 is a protein involved in allegraft differentiat rejection. In experimental autoimmune encephalomyelitis (EAE), neuritis(EAN) and ion/developm uveitis (EAU) it is produced by macrophages and microglia cells.	differentiat ion/developm ent
any2_2b19	any2_2b19 Originates from TXBP151 mRNA by alternative splicing	It is ubiquitously expressed. The mRNA is also subject to alternative polyadenylation. Overexpression of TXBPLS1 in NIH3T3 cells causes inhibition of apoptosis induced by tumour necrosis factor (TNF). It binds to A2O, which is also an inhibitor of cell death by a yet unknown mechanism.	differentiat ion/developm ent
amy2_7j5	amy2_7j5 similarity to Tspyll testis- specific Y-encoded-like protein of hus musculus	stis- TSPY genes are arranged in clusters on the Y chromosome of many mammalian species. differentiat protein of TSPY is believed to function in early spermatogenesis and is a candidate for GBY, ion/developm the putative gonadoblastoma-inducing gene on the Y. The TSPY family forms part of ent a superfamily, TTSN, with autosomal representatives, highly conserved in mammals and beyond.	differentiat ion/developm ent

Group Intracellular Transport and Trafficking

CloneID	Homology	Punotison	dronb
amy2_14b5	shows blk identity to the human TYL protein and 46% identity to the human Tic protein	of Saccharomyces cerevisiae, which takes new protein shows also significant the control of Golgi structure an matty conneced in the control of Golgi structure and the control of Golgi structure and the control of the	intracellula r transport and
amy2_2033	amy2_2013 high similarity to murine synaptotagmin 3.	The novel protein contains two C2 domains. The C2 domain is thought to be involved intracellula in C2 domains. The C2 domains are essential for r transport (a(2+)-regulated exocytosis of neurosecretory vesicles	intracellula intracellula r transport and trafficing
fkd2_3k1	very similar to rat testicular dynamin	Dynamin is a microtubule-associated force-producing protein, which is involved in intracellula the production of microtubule bundles and which is able to bind and hydrolyze GTP r transport and provides the motor for vesicular transport during endocytosis. The protein is and ubiquitously expressed, but in brain and testis above average.	intracellula r transport and trafficing
mel2_7g14	mel2_7g14 Similarity to the dor (deep orange) protein of drosophila melanogaster.	tain pep3 of mechanisms. The rt/targeting is	intracellula r transport and

Group Melanoma derived

_:: ₂			
Group	melanoma	melanoma	
Punotton	The novel protein contains a leucin zipper. No informative BLAST results: No needictive procette, near on trop motife.	Transcpripts can be found in almost any tissue, but are most abundant in kidney and retina.	INO INTOTORILVE BIANT TESTITOS NO DIRECTION DIRECTION DIRECTION
Homology	similarity to integrin I of Saccharomyces cerevisiae	without similarity to known proteins	
orreps	mele_lejl	mel2_7k39	

Group Metabolism

		Winction	dnozb
Sil	amy2_2c22 similarity to the L-acyl-glycerol- It contains	'yl-glycerol- It contains one leucine zipper. The protein is belived to play a role in fatty metabolism	metabolism
E E	mais. placenta and foreskin.	foreskin.	
Şį	aspartate	The L-isoaspartyl methyltransferase (Pimt), as an example, is a highly conserved metabolism	metabolism
9	methyltransferases.	enzyme utilising S-adenosylmethionine (AdoMet) to methylate aspartate residues of	
ļ	proteins dam	protains damaged by age-related isomerisation and deamidation.	

Group Nucleic acid management

Cloneid	Homology		
DKFZph		The second of th	droge
amy2_lln4	similarity to RADLA of		nucleic acid
	Schizosaccharomyces pombe and	DNA	annagement.
	YLR3å3w of Saccharomyces cerevisia.		, , , , , , , , , , , , , , , , , , ,
amy2_lil	similarity to the murine hemin-	The hemin-sensitive initiation factor 2 is expressed predominantly in liver.	nine nietnie
	sensitive initiation factor 2.	ologue	management
amy2_2912	similarity to NVL-2 of Rattus	. —	nucleic acid
	norvegicus.		management
		_	
	-00	activity-dependent manner. The new protein exhibists elevated expression in brain	٠.
FLAT TRAS		and testis.	
ו מניהם	Dre_roter night CSimilarity to glutamytRNA		nucleic acid
	A TURNING GIBLELOSE SCOULTE A	ut pund	management
	or the hyperthermophilic bacterium	numerous ATP- or GTP-binding protains, such as ATP synthase alpha and beta	
	Aquitex abolicus.	Subunits, Myosin heavy chains, Kinesin heavy chains and kinesin-like proteins.	
		Dynamins and dynamin-like proteins, several kinases, DNA and RNA helicases, GTP-	
		binding elongation factors and the Ras family of GTP-binding proteins. The protein	
		seems to be expressed ubiquitously.	
tes3_loil6	tes3_lOils similarity to human ZK1.	The ZK1 gene is one of early response genes by exposure to ionizing radiation, and nucleic acid	nucleic acid
		plays a role in radiation-induced apoptotic cell death on hematopoietic cells. The management	management
		noval protein contains 18 zinc finger domains, a RGD call attachment and a ATP GTP	
		A domain.	
tes3_31a10	tes3_3lald similarity to histone Hi of	Histone Hi variants are known to act as specific regulators of genes via the	nucleic acid
			management

Group Signal transduction

CloneID DKFZp.	Homology	Section Sectio	Group
amy2_10h17	weak similarity to murine hacl	ains a Zinc finger motif of the CBHC4 type (RING finger). is involved in mediating protein-protein interactions. RING-finger are: mammalian V(D) U recombination activating rpt-l, human rfp, human 52 Kd Ro/SS-A protein and others. ger proteins contains a number of oncogenes. For example ription factor, BRCAl, the mammalian cbl- and bmi-l proto-	signal transduction
amy2_10p7	similarity to Na+/Ca2+ exchange proteins	ort of (a2+ from the sarcoplasm into the sarcoplasmic reticulum is an process in the initiation of muscle relaxation. n. the novel protain contains a PROSITE multicopper oxidase signature. r oxidases are enzymes that possess three spectroscopically different ters.	signal transduction
amy2_b2d7	a so far unknown alternative spliced form of disks large homolog DLG2.	It seems to be predominantly expressed in the retina, germ cells and brain. It contains a SH3-domain and a guanylate kinase domain. These conserved regions are the shared among members of the discs-large family of proteins that include human p55, a membrane protein expressed in erythrocytes, rat PSD-95/SAP90, a synapse protein expressed in brain, Drosophila dIg-A, a septate junction protein expressed in proteins, and human and mouse 20-1 and canine 20-2, two tight junction proteins. The Homologue of Drosophila, dIg-A, acts as a tumor suppressor. All members of this family may be involved in signal transduction.	signal transduction
amy2_2fl&	TI 25 1	ial in the assembly, mplex of the rat to be restricted to	signal transduction
tes3_bbc22	Partial similarity to mouse PC326	e known as regulatory propeller like eraction. The new an essential	signal transduction
tes3_bld2l	Contains the full coding sequence of the human Nedd-4-like ubiquitin-protein ligase.	Www domains. The WW/rsp5/WWP domain has been shown r proline-motifs, and thus resembles somewhat SH3 clated with other domains typical for proteins in There is also a ubiquitin-protein ligase activity ad to play an important role in protein-degradation	signal transduction
tes3_29f24	Similarity to murine netla.	The closely related mNETL activates signalling pathways in addition to those sidirectly controlled by activated RhoA. The novel protein is expressed ubiquitously.	signal transduction
tes3_31,50	contains a Protein phosphatase 2C. motif.	The novel protein shares 95% identity withthe rat protein phosphatase 2C and is sexpressed ubsquitcusly. PPEC is a structurally diversified protein phosphatase family with a wide range of functions in cellular signal transduction. The transcription of the PPEC delta gene was activated in response to stress, like alcohol or UV irridation. PPEC plays a role in cell cycle control.	signal transduction
tes3_5k22	similarity to human paraneoplastic neuronal antigen MAI	neurological ain, but ESTs	signal transduction

Group Testis derived

	_	_	,								
Group	testis derived	testis	testis	testis	testis. derived	testis derived	testis	testis	testis derived	testis derived testis	derived testis derived
The state of the s	ially polyadenylated and the novel protein is u results: No predictive prosite, nfam or KOP mo	No informative BLAST results: No predictive prosite, pfam or SCOP motife.	The EST-distribution signifies an ubiquitous expression pattern. No informative BLAST results: No predictive prosite, pfam or KCOP motife.	The MRNA is transcribed ubiquitously. No informative BLAST results: No predictive proxite, ofam or CCOD motife.	Neurofilaments are the intermediate filaments specific to nervous tissue. They are probably essential to the tensile strength of the neuron, as well as to transport of molacules and organellas within the axon. Until now, ESTs of the novel mRNA could only be isolated from testes, germ cells and uterus.	tif. be		No informative BLAST results: No predictive prosite, pfam or SCOP motife.	No informative BLAST results; No predictive prosite, pfam or SCOP motife.	No informative BLAST results: No predictive prosite, pfam or SCOP motife. No informative BLAST results: No predictive presite, pfam or SCOP motife	informative BLAST results; No predictive prosite, pfam or SCOP motife.
Homology	res3_bOn10 without similarity to known proteins.	Tes3_llel7 without similarity to known proteins.	Tes3_b2db6 without similarity to known proteins	without similarity to known proteins.	Tes3_15n14 weak similarity to the neurofilament triplet M protein of the rat.	without similarity to known proteins	without similarity to known proteins	es3_21k14 without similarity to known proteins	Weak similarity to RCC1-like 6 exchanging factor RL6, UVRb (UVB-resistance protein) of Arabidopsis thaliana and to the murine retinitis plamentosa 6TPase requiator.	Similarity to the F-box protein FBLE of the rat. Without similarity to known	
CloneID DKFZpg	Tes3_bonjo	Tes3_llel7	Tes3_12dl&	Tes3_1417	Tes3_15n14	Tesa_16p3	Tesa_L9pl2	Tes3_21k14	Tesa_22ill	Tes3_22124 tes3_26q3	tes3_30pb

Group Transmembrane proteins

CloneID	l. Homology	Ogozo	Group
amy2_11d2	my2_lld2 Without similarity to known proteins	The noval protein contains 2 transmembrane regions. No informative Blast results: no predictive prosite, pfam or scope motife.	transmembran e proteins
amye_lebol7	smy2_J2Jo17 without similarity to known proteins.	The novel protein contains 1 transmembrane region. No informative BLAST results: No predictive prosite, pfam or SCOP motife.	Transmembran e proteins
amy2_114	myz_lily Similarity to the human 1(3)mbt protein homolog.	Nutations of the Drosophila 1(3)mbt gene lead to malignant brain tumors. The protein contains one transmembrane domain No informative BLAST results; No predictive prosite, pfam or SCOP motife	Transmembran e proteins
amy2_24c6	amy2_24c6 without similarity to known proteins	The novel protein contains 1 transmembrane region. No informative BLAST results; No predictive prosite, pfam or SCOP motife.	Transmembran e proteins

fbr2_7844	_	The novel protein contains 1 transmembrane region and a Cytochrome c family heme- Transmembran	Transmembran
	proceins.	binding site. No informative BLAST results: No predictive prosite, pfam or SCOP motife.	e proteins
tes3_lbal7	tes3_11al7 without similarity to known proteins	The novel protein contains 2 transmembrane regions and one leucine zipper. The protein is ubiquitously expressed with higher abundance in stomach, brain and testis.	Transmembran e proteins
+063 17:31		No informative BLAST results: No predictive prosite, pfam or SCOP motife.	
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	cess_witte without stailBrity to known proteins	The novel protein contains 2 transmembrane regions. ESTs can be found in testis. Transmembran retina and brain.	Transmembran profesion
TANE COURT		No informative BLAST results: No predictive prosite, pfam or SCOP motife.	
37 CO T	cess_cours attended started to known proteins	The novel protein contains 1 transmembrane region and two leucine zippers.	Transmembran
tes3_7nl2.	tes3_7nl2. Without similarity to known	The novel protein contains I transmembrane domain	e proteins Transmembran
+nc3 0,11		No informative BLAST results: No predictive prosite, plam or SCOP motife.	e proteins
	without Similarity to known	ain contains & transmembrane region. The only EST described so far	Transmembran
	6119001	No information Black note: No supplication security of the second of the	e proteins
		ing the course ocasi issuits, no predictive prosite, prom or stop motifie	

Group Transcription factors

				
Group	transcriptio n factors	transcriptio n factors	transcriptio n factors	Transcriptio n factors
And the second of the second o	similarity to the homeotic protein Homoeobox genes are known to play important roles in developmental processes. In transcription emx2 of man, mouse and zebra fish zebrafish emx2 mRNAs are found in the dorsal talencephalon, parts of the next mRNAs are found in the human homologue Emx2 appears to be already attented to the gene "empty diencephalon and the otocyst. The human homologue Emx2 appears to be already carefully spiracles of Drosophila expressed in a.5 day embryos. It is also expressed in the presumptive cerebral cortex, olfactory blactory blactory blactory blactory blactory blactory blactory blactory plactory in earlier stages and olfactory epithelia later in development. Mutants of the D. melanogaster gene "mempty spiracles" display spiracles devoid of filzkorper, no antenna and an open head.	I-kappa-B-related protein interacts with transcription factors and BRCA1 has a function in DNA damage response. I-kappa-B-alpha mutations contribute to constitutive NY-kappaB activity in cultured and primary NRS (Hodgkin/Reed-Sternberg) cells and are therefore involved in the pathogenesis of Hodgkin's disease (NP) patients	The novel protein is ubiquitously expressed. YDL153c is involved in transcriptional silencing.	Giantin is discussed as an autoantigen in rheumatoid arthritis. The novel protein Transcriptio contains a leucine zipper and a putative Helix-loop-helix DNA-binding domain. Therefore it might be a novel transkription factor. Nost EST hits are from tastis
Ношо Лоду	similarity to the homeotic protein emx2 of man, mouse and zebra fish as wall as to the gene "empty spiracles" of brosophila melangaster.	partial identity to I-kappa-B- related protein and to BKCAI.	any2_2f22 similarity to YDL153c of Saccharomyces cerevisia	tes3_bbnl4 similarity to human giantin.
CloneID DKFZp		amye_bche	amy2_2f22	tes3_lånl4

DKFZphamy2_10h17

group: signal transduction

10 DKFZphamy2_10hl? encodes a novel 180 amino acid protein which shows weak similarity to murine hacl.

The novel protein contains a Zinc finger motif of the C3HC4 type (RING finger). The RING-finger domain is involved in mediating protein-protein interactions. Proteins containing a RING-finger are: mammalian V(D)J recombination activating protein (RAGL); mouse rpt-L; human rfp; human 52 Kd Ro/SS-A protein and others. The family of RING finger proteins contains a number of oncogenes. For example PML; a probable transcription factor; BRCAL; the mammalian cbl- and bmi-L proto-oncogenes.

The new protein can find application in modulating proteinprotein-interaction and in studying the expression profile of amygdala-specific genes.

25

5

weak similarity to hacl (Mus musculus)

Sequenced by LMU

30

35

Locus: unknown

Insert length: 835 bp
Poly A stretch at pos. 751, polyadenylation signal at pos. 729

L CACAGAGATC ATTGTCAACC AGGCCTGTGG GGGGGACATG CCTGCCTTGG 51 AAGGGGCACC CCATACCCCG CCACTGCCAC GGCGGCCCCG TAAGGGAAGC LOL TCGGAGCTGG GCTTTCCCCG CGTGGCCCCA GAGGATGAGG TCATTGTGAA 40 151 TCAGTACGTG ATTCGGCCTG GCCCCTCGGC CTCGGCGGCT TCTTCGGCGG 201 CGGCAGGGA GCCCCTGGAG TGCCCCACCT GTGGGCACTC CTACAATGTC 251 ACCCAGCGGA GGCCCCGCGT GCTGTCCTGC CTGCACTCTG TGTGTGAGCA 3D1 GTGCCTGCAG ATTCTCTACG AGTCCTGCCC CAAGTACAAG TTCATCTCCT 351 GCCCCACCTG CCGCCGTGAG ACTGTGCTCT TCACCGACTA CGGCCTGGCC
401 GCGCTGGCTG TCAACACGTC CATCCTGAGC CGCCTGCCGC CTGAGGCGCT 45 451 GACGGCCCCA TCCGGGGGTC AGTGGGGGGC TGAGCCCGAG GGCAGCTGCT 501 ACCAGACCTT CCGGCAGTAC TGTGGGGCCG CGTGCACCTG CCACGTGCGG 551 AACCCACTGT CCGCCTGCTC CATCATGTAG TAGCGCCTGC CTGCCCGCCA LSI GCCGCCCGCT GACCCTTCCT TCCCCACCAT GGCTTCCGGC CCCACCCCGA 50 701 GTGGCATTGT CGCTGCAGCC AACTTTGCCA TTAAAACTCT TTGCCAAAGT BOL AAAAAAAAA AAAAGAAAAA AAAAAAAAA AAAAG

55

BLAST Results

No BLAST result

5

Medline entries

No Medline entry

10

Peptide information for frame 2

ORF from 38 bp to 577 bp; peptide length: 180
15 Category: similarity to unknown protein
Classification: Cellular transport and traffic
Prosite motifs: PRENYLATION (177-180)
ZINC_FINGER_C3HC4 (81-90)

20

1 MPALEGAPHT PPLPRRPRKG SSELGFPRVA PEDEVIVNQY VIRPGPSASA
51 ASSAAAGEPL ECPTCGHSYN VTQRRPRVLS CLHSVCEQCL QILYESCPKY
101 KFISCPTCRR ETVLFTDYGL AALAVNTSIL SRLPPEALTA PSGGQWGAEP
151 EGSCYQTFRQ YCGAACTCHV RNPLSACSIM

25

BLASTP hits

30 No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_10h17, frame 2

No Alert BLASTP hits found

35

Pedant information for DKFZphamy2_10h17, frame 2

Report for DKFZphamy2_10h17-2

40

50

ELENGTHJ 180 EMUJ 19400-27 Epij 7-95

45 EHOMOLI TREMBL:ACOD7727_7 gene: "F&K7.7"; Arabidopsis thaliana chromosome I BAC F&K7 sequence; complete sequence. 3e-06

EBLOCKZI BLOOA39C
EBLOCKZI BLOOA39C

EBLOCKSI BLOO518 Zinc finger, C3HC4 type, proteins
EPROSITEI PRENYLATION 1

EPROSITE ZINC_FINGER_C3HC4 1

IPFAMD Zinc finger, C3HC4 type (RING finger)

55 [KW] Alpha_Beta LOW_COMPLEXITY 5.56 %

	W	O 01/98454	PCT/IB01/02050			
	SEQ SEG PRD	MPALEGAPHTPPLPRRPRKGSSELGFPRVAPE	,			
5	SEQ SEG PRD	ECPTCGHSYNVT@RRPRVLSCLHSVCE@CL@I				
10	SEQ SEG PRD	AALAVNTSILSRLPPEALTAPSGGQWGAEPEG				
15	Prosite for DKFZphamy2_lDhl7.2					
	0029 0029					
20		Pfam for DKFZph	hamy2_10h17.2			
25	нмм_	NAME Zinc finger, C3HC4 type (RI	ING finger)			
30	HMM *CPICFcTF@1DyPWPFdePmM1PCgHsFCypCIrrWC CP C Y+ +P+ L C+HS C+ C+ ++ OC Query L2 CPTCGHSYNVT@RRPRVLSCLHSVCE@CL- QILYESCPKYKFISC 105					
35	HMM Quer	PmC* P C y 106 PTC 108				

DKFZphamy2_10p7

5 group: signal transduction

DKFZphamy2_10p7 encodes a novel 1615 amino acid protein with similarity to Na+/Ca2+ exchange proteins.

10 The Transport of Ca2+ from the sarcoplasm into the sarcoplasmic reticulum is an essential process in the initiation of muscle relaxation.

In addition, the novel protein contains a PROSITE multicopper oxidase signature. Multicopper oxidases are enzymes that possess

15 three spectroscopically different copper centers.

The new protein can find application in modulation of NA+/Ca2+- exchange and voltage-dependend processes.

similarity to Na+/Ca2+ exchange proteins

ATG in frame 3 is first in clone.

25 Sequenced by LMU

Locus: unknown

Insert length: 5236 bp
30 Poly A stretch at pos. 5216, no polyadenylation signal found

L CGGACGCGTG GGCGGACGCG TGGGCCCTGT ATACCTGTGC CACTTTGTGC 51 CTTAAGGAAC AAGCTTGCTC AGCGTTTTCA TTTTTCAGTG CTTCTGAGGG LOL TCCCCAGTGT TTCTGGATGA CATCATGGAT CAGCCCAGCT GTCAACAATT 35 151 CAGACTTCTG GACCTACAGG AAAAACATGA CCAGGGTAGC ATCTCTTTTT 201 AGTGGTCAGG CTGTGGCTGG GAGTGACTAT GAGCCTGTGA CAAGGCAATG 251 GGCCATAATG CAGGAAGGTG ATGAATTCGC AAATCTCACA GTGTCTATTC 301 TTCCTGATGA TTTCCCAGAG ATGGATGAGA GTTTTCTAAT TTCTCTCCTT 40 351 GAAGTTCACC TCATGAACAT TTCAGCCAGT TTGAAAAATC AGCCAACCAT 4Dl AGGACAGCCA AATATTTCTA CAGTTGTCAT AGCACTAAAT GGTGATGCCT 451 TTGGAGTGTT TGTGATCTAC AGTATTAGTC CCAATACTTC CGAAGATGGC 5DL TTATTTGTTG AAGTTCAGGA GCAGCCCCAA ACCTTGGTGG AGCTGATGAT 551 ACACAGGACA GGGGGCAGCT TAGGTCAAGT GGCAGTCGAA TGGCGTGTTG LOL TTGGTGGAAC AGCTACTGAA GGTTTAGATT TTATAGGTGC TGGAGAGATT 45 651 CTGACCTTTG CTGAAGGTGA AACCAAAAAG ACAGTCATTT TAACCATCTT 7DL GGATGACTCT GAACCAGAGG ATGACGAAAG TATCATAGTT AGTTTGGTGT 751 ACACTGAAGG TGGAAGTAGA ATTTTGCCAA GCTCCGACAC TGTTAGAGTG BDL AACATTTTGG CCAATGACAA TGTGGCAGGA ATTGTTAGCT TTCAGACAGC 851 TTCCAGATCT GTCATAGGTC ATGAAGGAGA AATTTTACAA TTCCATGTGA 50 TAAGAACTTT CCCTGGTCGA GGAAATGTTA CTGTTAACTG GAAAATTATT 951 GGGCAAAATC TAGAACTCAA TTTTGCTAAC TTTAGCGGAC AACTTTTCTT LOOL TCCTGAGGGG TCGTTGAATA CAACATTGTT TGTGCATTTG TTGGATGACA 1051 ACATTCCTGA GGAGAAAGAA GTATACCAAG TCATTCTGTA TGATGTCAGG 55 1251 ATTCAGCTTT TCATCAACAG AGAATTTGGA TCTCTAGGAG CTATCAATGT

LED CACATATACC ACGGTTCCTG GAATGCTGAG TCTGAAGAAC CAAACAGTAG 1351 GAAACCTAGC AGAGCCAGAA GTTGATTTTG TCCCTATCAT TGGCTTTCTG 1401 ATTTTAGAAG AAGGGGAAAC AGCAGCAGCC ATCAACATTA CCATTCTTGA 1451 GGATGATGTA CCAGAGCTAG AAGAATATTT CCTGGTGAAT TTAACTTACG 5 1501 TTGGACTTAC CATGGCTGCT TCAACTTCAT TTCCTCCCAG ACTAGATTCA 1551 GAAGGTTTGA CTGCACAAGT TATTATTGAT GCCAATGATG GGGCCCGAGG ILDI TGTAATTGAA TGGCAACAAA GCAGGTTTGA AGTAAATGAA ACCCATGGAA 1651 GTTTAACATT GGTAGCCCAG AGGAGCAGAG AACCTCTTGG CCATGTTTCC 10 15 20 25 45 50 3801 GCTGATGGAA GTATTAGTGA TATATTTCCT ACCTCAGGAG TGATTTTATT 3851 TACTGAAGGC CAGGTACTGT CAACAATCAC TCTAACTATT CTTGCTGATA 390% ATATACCAGA GTTATCAGAG GTTGTGATTG TAACCCTCAC CCGTATCACC 395% ACAGAAGGG TTGAGGACTC ATACAAGGT GCTACTATTG ATCAGGACAG 4001 AAGCAAGTCT GTTATAACAA CTTTGCCCAA TGACTCACCT TTTGGCTTGG 55 4051 TGGGCTGGCG TGCTGCGTCT GTCTTCATTA GAGTAGCAGA GCCTAAAGAA 4101 AACACCACCA CTCTTCAGTT ACAAATAGCT CGAGATAAAG GACTACTTGG 4151 GGATATTGCC ATTCACTTGA GAGCTCAACC CAATTTCTTA CTGCATGTCG

WO 01/98454
PCT/IB01/02050
H2Dl ATAATCAAGC TACTGAGAAT GAAGATTATG TATTGCAAGA AACAATAATA
H25l ATAATGAAAG AAAACATAAA AGAAGCTCAT GCCGAAGTTT CCATTTTGCC
H3Dl GGATGACCTT CCTGAATTGG AGGAAGGATT TATTGTCACT ATCACTGAGG

4351 TGAACCTGGT GAACTCTGAC TTCTCTACAG GACAGCCAAG TGTGCGGAGG
5 4401 CCCGGAATGG AAATAGCTGA GATAATGATA GAAGAAAATG ACGATCCCAG
4451 AGGAATTTTT ATGTTTCATG TTACTAGAGG CGCTGGGGAA GTTATTACTG
4501 CCTATGAGGT GCCTCCACCC TTGAACGTTC TTCAAGTTCC TGTAGTCCGG

4501 CCTATGAGGT GCCTCCACCC TIGAACGTTC TICAAGTTCC TGTAGTCCGG
4551 CTGGCTGGAA GCTTTGGGGC AGTAAATGTT TATTGGAAAG CATCACCAGA
4601 CAGTGCTGGC CTGGAAGACT TTAAACCATC TCATGGGATT CTTGAATTTG

10 4651 CAGATAAACA GGTTACTGCA ATGATAGAAA TCACCATAAT TGATGATGCT 4701 GAATTTGAAT TGACAGAGAC GTTCAATATT TCCTTGATCA GTGTTGCTGG 4751 AGGTGGCAGA CTTGGTGATG ATGTTGTGGT AACTGTTGTT ATTCCACAAA

4803 ATGATTCTCC ATTTGGAGTA TTTGGATTTG AAGAAAAGAC TGTAAGTTAA
4853 ACATATCAGG GGAAAGCCTT GTTTCAGGCT AGCGTTTCAT GTAATTTTGA

15 4901 GTAGAAAGTG TCTCACATTT TTGTTTTGGA AGTCTTGGCC AGGCATGGTG
4951 GCTCATGCCA GTAATCCCAG CACTTTGGGA GGCCGCAGCG GGCAGATCAC
5001 GAGGTCAGGA GATTGACACC ATCCTGGCCA ATATGGTTGA ATTCCCGTCT

5001 GAGGTCAGGA GATTGACACC ATCCTGGCCA ATATGGTTGA ATTCCCGTCT 5051 CTACTGAAAG TACAAAAATT AGCTGGGCGT GGTGGCACAT GCCTGTATTC 5101 CCAGATACTT GGGAGGCTGA GGCAGGAGAC TCGCTTGAAC CCAGGAGGCA

20 5151 GAGGTTGCAG TGAGCTGAGA TCACGCCATT GCACTCCAGC CTGGCGACAT 5201 AGAGAGACTC CATCTCAAAA AAAAAAAAA AAAAAG

BLAST Results

No BLAST result

30 Medline entries

No Medline entry

35

25

Peptide information for frame 3

40 ORF from O bp to 4847 bp: peptide length: 1616 Category: putative protein Classification: Cell signaling/communication Prosite motifs: MULTICOPPER_OXIDASEL (151-171)

1 DAWADAWALY TCATLCLKEQ ACSAFSFFSA SEGPQCFWMT SWISPAVNNS
51 DFWTYRKNMT RVASLFSGQA VAGSDYEPVT RQWAIMQEGD EFANLTVSIL
101 PDDFPEMDES FLISLLEVHL MNISASLKNQ PTIGQPNIST VVIALNGDAF
151 GVFVIYSISP NTSEDGLFVE VQEQPQTLVE LMIHRTGGSL GQVAVEWRVV
251 TEGGSRILPS SDTVRVNILA NDNVAGIVSF QTASRSVIGH EGEILQFHVI
301 RTFPGRGNVT VNWKIIGQNL ELNFANFSGQ LFFPEGSLNT TLFVHLLDDN
351 IPEEKEVYQV ILYDVRTQGV PPAGIALLDA QGYAAVLTVE ASDEPHGVLN
401 FALSSRFVLL QEANITIQLF INREFGSLGA INVTYTTVPG MLSLKNQTVG
551 NLAEPEVDFV PIIGFLILEE GETAAAINIT ILEDDVPELE EYFLVNLTYV
501 GLTMAASTSF PPRLDSEGLT AQVIIDANDG ARGVIEWQQS RFEVNETHGS
551 LTLVAQRSRE PLGHVSLFVY AQNLEAQVGL DYIFTPMILH FADGERYKNV
601 NIMILDDDIP EGDEKFQLIL TNPSPGLELG KNTIALIIVL ANDDGPGVLS

WO 01/98454 PCT/IB01/02050 651 FNNSEHFFLR EPTALYVQES VAVLYIVREP AGGLFGTVTV QFIVTEVNSS 701 NESKDLTPSK GYIVLEEGVR FKALQISAIL DTEPEMDEYF VCTLFNPTGG 751 ARLGVHVQTL ITVLQNQAPL GLFSISAVEN RATSIDIEEA NRTVYLNVSR ADD TNGIDLAVSV QUETVSETAF GMRGMDVVFS VFQSFLDESA SGUCFFTLEN 5 B51 LIYGIMLRKS SVTVYRWQGI FIPVEDLNIE NPKTCEAFNI GFSPYFVITH 901 EERNEEKPSL NSVFTFTSGF KLFLVQTIII LESSQVRYFT SDSQDYLIIA 951 SQRDDSELTQ VFRWNGGSFV LHQKLPVRGV LTVALFNKGG SVFLAISQAN IDDI ARLNSLLFRW SGSGFINFQE VPVSGTTEVE ALSSANDIYL IFAKNVFLGD 1051 QNSIDIFIWE MGQSSFRYFQ SVDFAAVNRI HSFTPASGIA HILLIGQDMS 10 1101 ALYCUNSERN QFSFVLEVPS AYDVASVTVK SLNSSKNLIA LVGAHSHIYE 1151 LAYISSHSDF IPSSGELIFE PGEREATIAV NILDDTVPEK EESFKVQLKN 1501 PKGGAEIGIN DZVTITILZN DDAYGIVAFA QNZLYKQVEE MEQDZLVTLN 1251 VERLKGTYGR ITIAWEADGS ISDIFPTSGV ILFTEGQVLS TITLTILADN 1301 IPELSEVVIV TLTRITTEGV EDSYKGATID QDRSKSVITT LPNDSPFGLV 1351 GWRAASVFIR VAEPKENTTT LQLQIARDKG LLGDIAIHLR AQPNFLLHVD 15 1401 NGATENEDYV LGETIIIMKE NIKEAHAEVS ILPDDLPELE EGFIVTITEV 1451 NLVNSDFSTG QPSVRRPGME IAEIMIEEND DPRGIFMFHV TRGAGEVITA 1501 YEVPPPLNVL QVPVVRLAGS FGAVNVYWKA SPDSAGLEDF KPSHGILEFA 1551 DKQVTAMIEI TIIDDAEFEL TETFNISLIS VAGGGRLGDD VVVTVVIPQN

BLASTP hits

25

20

No BLASTP hits available

1601 DSPFGVFGFE EKTVS

Alert BLASTP hits for DKFZphamy2_10p7, frame 3

30 TREMBL:AF055084_l gene: "VLGRl"; product: "very large G-protein
coupled
receptor-l"; Homo sapiens very large G-protein coupled receptorl
(VLGRL) mRNA; complete cds.; N = 3; Score = 284; P = 1.2e-33

TREMBL:DMAF9897_1 gene: "Calx"; product: "CALX"; Drosophila melanogaster 3Na(+)-1Ca(2+) exchanger (Calx) mRNA, complete cds., N =

l. Score = 178. P = 3.3e-09

40

35

>TREMBL:AFD55D84_1 gene: "VLGR1"; product: "very large G-protein coupled

receptor-l": Homo sapiens very large G-protein coupled 45 receptor-l (VLGR1)

mRNA: complete cds:

Length = 1,967

HSPs:

50

Score = 284 (42.6 bits), Expect = 1.2e-33, Sum P(3) = 1.2e-33 Identities = 192/738 (26%), Positives = 314/738 (42%)

Query: 57
55 SGQAVAGSDYEPVTRQWAIMQEGDEFANLTVSILPDDFPEMDESFLISLLEVHLMNISAS 12b

S + G DY + Q G + + +SI+ D+ E +E +E+

102 SZASPGGVDYI-LHGSTVTFQHGQNLSFINISIIDDNESEFEEP----Sbict: IEILLTGATGG 155

Querv: 752

5 LKNQPTIGQPNISTVVIALNGDAFGVFVIYSISPNTSEDGLFVEVQEQPQTLV-ELMIHR 185 +G+ +S ++IA + FGV N Z+

T++ L++ R

Sbjct: 15b A----VLGRHLVSRIIIAKSDSPFGVIRFL----NgSK----ISIANPNSIMILSLVLER 203

10

186 TGGSLGQVAVEWRVVGGTATEGL----DFIG-AGEILTFAEGETK-Query: RTVILTXXXXXXX 238

TGG LG++ V W VG + E L D F EGE

+T+ILT

15 Sbjct: 204

TGGLLGEIQVNWETVGPNZQEALLPQNRDIADPVZGLFYFGEGEGGVRTIILTIYPHEEI 263

239 XXXXXXXXLVYTEGGSRILPSSDTVRVNILANDNVAGIVSF--QTASRSVIGH----EG 292

20 EG

+6 +++ + V + I+ G+V F +T S+

Sbjct: 264 EVEETFIIKLHLVKGEAKLDSRAKDVTLTIQEFGDPNGVVQFAPETLSKKTYSEPLALEG 323

25 293 EILQFHVIRTFPGR-GNVTVNWKIIGQ-NLELNFANFSGQLFFPEGSLNTTLFVHLLDDN 350 +L +R G G + V W++

+ ++ +F + SG +6

VHLL D 324 Sbjct:

30 PLLITFFVRRVKGTFGEIMVYWELSSEFDITEDFLSTSGFFTIADGESEASFDVHLLPDE 383

Query: 351 IPEEKEVYQVILYDVRTQGVPPAGIALLDAQGYAAVLTVEASDEPHGVLNFAL-SSRFVL 409 +PE +E Y + L V G A LD + · +V A+D+PHGV

35 FAL S R 384 VPEIEEDYVIQLVSVE-----GGAELDLEKSITWFSVYANDDPHGV--Sbjct: FALYSDRQSI 434

410 LQEANI--TIQLFINREFGSLGAINVTYTTVPGMLSLKNQT-Query:

VGNLAEPEVDFVPIIGFL 466 40

N+ +IQ+ IR G+ G+ VAE +

Sbjct: 435 LIGANLIRSIAINITRLAGTE GDVAVGLRISSDH---KEAPIVTENAERA------L 482

45 Query: 467

ILEEGETAAAINITILEDDVPELEEYFLVNLTYVGLTMAASTSFPPRLDSEGLTAQVIID 526 ++++G T + I L F+L VL +A V+

50 ° Sbjct: 483 VVKDGATYKVDVVPIKNQVFLSLGSNFTLQLVTVMLVGGRFYGMPTILQ-EAKSA-VLPV 540

527 ANDGARGVIEWQQSRFEV-NETHGSLTLVAQRSREPLGHVSLFV---Query: YAQNLEAQVGLDY 582

55 + ++ + F++ N T G+ ++ R R G +S+ YA. LE Sbict: 541 SEKAANSQVGFESTAFQLMNITAGTSHVMISR-RGTYGALSVAWTTGYAPGLEIPEFIVV 599

```
Query:
             583 -IFTPMI--
    LHFADGERYKNVNIMILDDDIPEGDEKFQLILTNPSPGLELGKNTIALIIV L39
                    TP + L F+ GE+ K V +
                                                   EFLL+
                                                                   G
 5
    Sbjct:
             600 GNMTPTLGSLSFSHGEQRKGVFLWTFPS--
    PGWPEAFVLHLSGVQSSAPGGAQLRSGFIV 657
             640 LANDDGPGVLSFN-
10
    NSEHFFLREPTALYVQESVAVLYIVREPAQGLFGTVTVQFIVTEVN 698
                  A + GV F+.+S + + E T + ++ V L+
    Sbjct:
             L58 -AEIEPMGVFQFSTSSRNIIVSEDTQM-IRLHVQRLF-----
    GFHSDLIKVSYQTTAG 708
15
    Query:
            L99 ZSNESKDLTP-SKGYIVLEEGVRFKALQISAILDTEPEMDEYFVCTL----
    ----FNP 747
                              G + ++
                                            +I+ID
                                                      E++E+F
    F+
20
    Sbjct:
             709
    SAKPLEDFEPV@NGELFF@KF@TEVDFEITIIND@LSEIEEFFYINLTSVEIRGL@KFDV 768
             748 TGGARLGVHVQT-LITVLQNQAPLGLFSISAVENR-ATSIDIE----
    EANRTVYLNVSRT 801
25
                     RL +
                              +IT+L N
                                         G+ IS E A ++D
    T+2 +1Y
             769 NWSPRLNLDFSVAVITILDNDDLAGM-
    Sbjct:
    DISFPETTVAVAVDTTLIPVETESTTYLSTSKT 827
30
             802 NGI 804
    Query:
    Sbjct:
             DEB ITT 850
     Score = 266 (39.9 bits), Expect = 4.0e-25, Sum P(3) = 4.0e-25
35
     Identities = 175/708 (24%), Positives = 306/708 (43%)
    Query:
             131
    PTIGQPNISTVVIALNGDAFGVFVIYSISPNTSEDGLFVEVQEQPQTLVELMIHRTGGSL 190
                 P IG +I ++I N +A G+
                                             P +
                                                        EV+E
40
    + R G+
    Sbjct:
              39 PEIGNISIVRIIIMKNDNAEGII---EFDPKYTA----FEVEEDVG-
    LIMIPVVRLHGTY 90
             191 GQVAVEWRVVGGTATEG-
    Query:
    LDFIGAGEILTFAEGETKKTVILTXXXXXXXXXXXXXXXXX 249
45
                 6 V ++
                             +A+ G +D+I G +TF G+
    Sbjct:
              91
    GYVTADFIZ@ZZZAZPGGVDYILHGZTVTF@HG@NLZFINIZIIDDNEZEFEEPIEILLT 150
50
            250 YTEGGSRILPSSDTVRVNILANDNVAGIVSFQTASRSVIGHEGE--
    Query:
    ILQFHVIRTFPGRG 307
                    GG+ +L
                                R+I+D+G++F
                                                   Z+
                                                       I +
                                                               IL
    RT
55
    Sbjct:
             151 GATGGA-
    VLGRHLVSRIIIAKSDSPFGVIRFLNQSKISIANPNSTMILSLVLERTGGLLG 209
```

WO 01/98454 PCT/IB01/02050 3DB NVTVNUKIIGQN-----LELN--FAN-FSGQLFFPEGSLNT-Querv: TLFVHLLDDNIPEEKEVY 358 + VNU+ +G N L N A+ SG +F EG T+++E +E + 5 Sbict: 570 EIQVNWETVGPNSQEALLPQNRDIADPVSGLFYFGEGEGGVRTIILTIYPHEEIEVEETF 2b9 359 QVILYDVRTQGVPPAGIALLDAQGYAAVLTVEASDEPHGVLNFA---Query: LSSRFV---LLQE 412 10 + L+ V+ G A LD++ LT++ +P+GV+ FA LZ + Sbjct: 270 IIKLHLVK-----GEAKLDSRAKDVTLTIQEFGDPNGVVQFAPETLSKKTYSEPLALE 322 15 413 ANITIQLFINREFGSLGAINVTYTTVPGMLSLKNQTVGNLAEPEVDFVPIIGFLILEEGE 472 + I F+ R G+ G I V + L ++ ++ E DF+ GF + +GE Sbict: 323 GPLLITFFVRRVKGTFGEINVVW-----EL22EF--DITE---20 DFLSTSGFFTIADGE 370 473 Query: TAAAINITILEDDVPELEEYFLVNLTYVGLTMAASTSFPPRLDSEGLTAQVIIDANDGAR 532 + A+ ++ +L D+VPE+EE +++ L 2 25 + AND Sbjct: 371 SEASFDVHLLPDEVPEIEEDYVIQLV-----SVEGGAELDLEKSITWFSVYANDDPH 422 Querv: 533 GVIEWQQSRFEV---NETHGSLTLVAQRSREPLGHVS--30 LFVYAQNLEAQVGLDYIFTPM 587 GV R + + +2 R G V+ L++E + . Sbjct: 423 GVFALYSDRQSILIGQNLIRSIQINITRLAGTFGDVAVGLRISSDHKEQPIVTENAERQL 482 35 588 ILHFADGERYKNVNIMILDDDI--PEGDE-KFQLILTNPSPGLELGKNTI--Query: -ALIIVLA 641 ++ DG YK V+++ + + + G QL+ G TI VI. 40 Sbict: 483 VVK--DGATYK-VDVVPIKNQVFLSLGSNFTLQLVTVMLVGGRFYGMPTILQEAKSAVLP 539 Querv: 642 NDDGPGVLSFNNSEHFFLREPTALYVQESVAVLYIVREPAQGLFGTVTVQFIV----TE L9L 45 NS+ F E TA + A ' +6 +6 ++V + V Ε 540 VSEKAA----NSQVGF--ESTAFQLMNITAGTSHVMISRRGTYGALSVAWTTGYAPGLE 592 50 Query: 697 VNSSNESKDLTPSKGYIVLEEGVRFKALQISAILDTEPEMDEYFVCTLFNPTGGARLGVH 756 ++TP+ G + G + K + + E FV L A G Sbjct: 593 IPEFIVVGNMTPTLGSLSFSHGEQRKGVFLWTF--55 PSPGWPEAFVLHLSGVQSSAPGGAQ 650

DLAVSVQWET B14

757 VQTLITVLQNQAPLGLFSISAVENRATSIDIEEANRTVYLNVSRTNGI--



V + + P+G+F +R +I + E + + L + V R G+ V ++T Sb ict: LSI LRSGFIVAEIE-PMGVFQFST-SSR--NIIVSEDT@MIRLHV@RLFGFHSDL-IKVSY@T 705 **5** . Querv: 815 VSETAFGMRGMDVVFS---VFQSFLDE 838 + +A + + V + FQ F E 706 TAGSAKPLEDFEPVQNGELFFQKFQTE 732 Sbict: 10 Score = 246 (36.9 bits), Expect = 4.le-32, Sum P(3) = 4.le-32 Identities = 92/338 (27%), Positives = 157/338 (46%) 511 PPRLDSEGLTAQVIIDANDGARGVIEW--QQSRFEVNETHGSLTLVAQRSREPLGHVSLF 568 15 PP + + + ++II ND A G+IE+ + + FEV E G + + G+V+ Sbjct: 38 PPEIGNISIV-RIIIMKNDNAEGIIEFDPKYTAFEVEEDVGLIMIPVVRLHGTYGYVTAD 96 20 Query: 569 VYAQNLEAQVG-LDYIFTPMILHFADGERYKNVNIMILDDDIPEGDEKFQLILTNPSPGL 627 +Q+AG+DYI+ F G+ +NI I+DD+ E +E +++LT + G 97 Sbict: 25 FIZQSSSASPGGVDYILHGSTVTFQHGQNLSFINISIIDDNESEFEEPIEILLTGATGGA 156 Query: P59 ELGKNTIALIIVLANDDGPGVLSFNNSEHFFLREPTALYVQESVAVLYIVREPAQGLFGT 687 GV+ F N LG++ ++ II+ +D 2 +L +V E 30 GL G Sbjct: 157 VLGRHLVSRIIIAKSDSPFGVIRFLNQSKISIANPN-----STMILSLVLERTGGLLGE 210 LAB VTVQFIVTEVNSSN----ESKDLT-PSKGYIVLEEGVR-Querv: 35 FKALQISAILDTEPEMDEYFV 741 + V + ZN +++D+ P G EG + + ++ Ε E++E F+ Sbjct: 577 IQVNWETVGPNSQEALLPQNRDIADPVSGLFYFGEGEGGVRTIILTIYPHEEIEVEETFI 270 40 742 CTLFNPTGGARLGVHVQTL-ITVLQNQAPLGL--FSISAVENRATSIDIE-EANRTVYLN 797 G A+L + + + T + +P G+ F+ 45 Sb.ict: IKLHLVKGEAKLDSRAKDVTLTIQEFGDPNGVVQFAPETLSKKTYSEPLALEGPLLITFF 330 798 VSRTNGIDLAVSVQWETVSETAFGMRGMDVVFSVFQSFLDESASGWCFFTL Query: 848 50 VRG + V ME ZE F + + FL S SG FFT+ 331 VRRVKGTFGEIMYVWELSSE------FDITEDFL--STSG--FFTI Sbjct: 366 Score = 24b (3b.9 bits), Expect = 1.9e-19, Sum P(3) = 1.9e-19 55 Identities = 87/303 (28%), Positives = 138/303 (45%) Query: 1162 PSSGELIFEPGEREA-TIAVNILDDTVPEKEESFKVQLKNPKGGAEIGIN-PLST STILLAST

WO 01/98454 PCT/1B01/02050 P ZG F GE TI + I E EE+F ++L KG A++ VT+TI Sbict: 536 PVSGLFYFGEGEGGVRTIILTIYPHEEIEVEETFIIKLHLVKGEAKLDSRAKDVTLTIQE 295 Query: 1220 NDDAYGIVAFAQNSL----YKQVEEMEQDSLVTLNVERLKGTYGRITIAWEADGSIS--- 1272 D G+V FA +L Y + +E L+T V R+KGT+G I + WE Sbict: 10 FGDPNGVVQFAPETLSKKTYSEPLALEGPLLITFFVRRVKGTFGEIMVYWELSSEFDITE 355 Query: 1273 DIFPTSGVILFTEGQVLSTITLTILADNIPELSEVVIVTLTRITTEGVEDSYKGATIDQD 1332 D TSG +G+ ++ + +L D +PE+ E ++ L ++ EG GA +D + 15 Sbjct: 35b DFLSTSGFFTIADGESEASFDVHLLPDEVPEIEEDYVIQL--VSVEG-----GAELDLE 407 Query: 1333 20 RSKSVITTLPNDSPFGLVGWRAASVFIRVAEPKENTTTLQLQIARDKGLLGDIAIHLRAQ 1392 + 2 + 2+ND P G+ I + +++Q+ I R G GD+A+ LR Sbjct: 408 KSITWFSVYANDDPHGVFALYSDRQSILIGQ--NLIRSIQINITRLAGTFGDVAVGLRIS 465 25 Query: 1393 PNFLLHVDNQ-ATENEDYVLQETIIIMKENIKEAHAEVSILPDDLPELEEGFIVTITEVN 1451 H + TEN E +++K+ VI F 30 Sbict: 466 SD---HKERPIVTENA----ERQLVVKDGATYKVDVVPIKNQVFLSLGSNFTLQLVTVM 517 Query: 1452 LVNSDFSTGQPSV 1464 LV F G P++ 35 Sbict: 518 LVGGRFY-GMPTI 529 Score = 246 (36.9 bits), Expect = 1.9e-19, Sum P(3) = 1.9e-19Identities = 89/334 (26%), Positives = 150/334 (44%) 40 Querv: 1159 DFIPSSGELIFEPGEREATIAVNILDDTVPEKEESFKVQLKNPKGGAEIGINDSVTITIL 1218 D+I + F+ G+ + I ++I+DD E EE ++ L GGA +G + Sbjct: 770

DYILHGSTVTFQHGQNLSFINISIIDDNESEFEEPIEILLTGATGGAVLGRHLVSRIIIA 169 45

Querv: 1219 SNDDAYGIVAFAQNSLYKQVEEMEQDSLVTLNVERLKGTYGRITIAWEADGSIS---- 1272 +D +G++ F S + +++L +ER G G I + WE

50 Z Sbict: 170 KSDSPFGVIRFLNQSKIS-IANPNSTMILSLVLERTGGLLGEIQVNWETVGPNSQEALLP 228

Query: 1273 ---DIF-PTSGVILFTEGQV-55 LSTITLTILADNIPELSEVVIVTLTRITTEGVEDSYKGA 1327 DI P SG+ F EG+ + TI LTI E+ E I+ L ΣŒ

WO 01/98454 PCT/IB01/02050 Sbjct: 229 QNRDIADPVSGLFYFGEGEGGVRTIILTIYPHEEIEVEETFIIKLHLVKGEAKLDS---- 284 Query: 1328 TIDQDRSKSVITTLPN-DSPFGLVGWRAASVFIRV-AEPK--5 ENTITLALAIARDKGLLG 1383 R+K V .T+ P G+V + +EP Ε R KG G Sbjct: 285 ----RAKDVTLTIQEFGDPNGVVQFAPETLSKKTYSEPLALEGPLLITFFVRRVKGTFG 339 10 Query: 1384 DIAIHLRAQPNFLLHVDNQATENEDYVLQETIIIMKENIKEAHAEVSILPDDLPELEEGF 1443 +I ++: F + ED++ +LPD++PE+EE + 15 Sbict: 340 EIMVYWELSSEFDI-----TEDFLSTSGFFTIADGESEASFDVHLLPDEVPEIEEDY 391 1444 IVTITEVNLVNSDFSTGQPSVRRPGMEIAEIMIEENDDPRGIFMFHVTR Query: 1492 20 + + I + NDDP G+F 392 VIQLVSVE-----GGAELDLEK---SITWFSVYANDDPHGVFALYSDR Sbjct: 431 Score = 237 (35.6 bits), Expect = 9.4e-34, Sum P(3) = 9.4e-34 25 Identities = 101/367 (27%), Positives = 165/367 (44%) SGRAVAGSDYEPVTRQWAIMREGDEFANLTVSILPDDFPEMDESFLISLLEVHLMNISAS 126 + G DY + Q G 30 Sbjct: 102 SSASPGGVDYI-LHGSTVTFQHGQNLSFINISIIDDNESEFEEP----IEILLTGATGG 155 Query: 153 35 LKNQPTIGQPNISTVVIALNGDAFGVFVIYSISPNTSEDGLFVEVQEQPQTLVELMIHRT 186 +G+ +Z ++IA + FGV N Z+ ++ L++ RT Sbjct: 15b A----VLGRHLVSRIIIAKSDSPFGVIRFL----NQSKISI---ANPNSTMILSLVLERT 204 40 187 GGSLGQVAVEWRVVGGTATEGL----DFIG-AGEILTFAEGETK-Query: PES XXXXXXXX 239 GG LG++ V W VG + E L D F EGE +T+ILT Sbict: 205 45 GGLLGEIQVNWETVGPNSQEALLPQNRDIADPVSGLFYFGEGEGGVRTIILTIYPHEEIE 264 240 XXXXXXXLVYTEGGZRILPSSDTVRVNILANDNVAGIVSF--QTASRSVIGH----EGE 293 +6 +++ V + IG+V F +T S+ 50 EG Sbict:

+F + SG

+G.

VEETFIIKLHLVKGEAKLDSRAKDVTLTIQEFGDPNGVVQFAPETLSKKTYSEPLALEGP 324

6 6 + V W++

294 ILQFHVIRTFPGR-GNVTVNWKIIGQ-

NLELNFANFSGQLFFPEGSLNTTLFVHLLDDNI 351

+R

+L

55

VHLL D +

Sbjct: 325

LLITFFVRRVKGTFGEIMVYWELSSEFDITEDFLSTSGFFTIADGESEASFDVHLLPDEV 384

Query: 352

5 PEEKEVYQVILYDVRTQGVPPAGIALLDAQGYAAVLTVEASDEPHGVLNFAL-SSRFVLL 410
PE +E Y + L V G A LD + +V A+D+PHGV FAL

Sbjct: 385 PEIEEDYVIQLVSVE-----GGAELDLEKSITWFSVYANDDPHGV--FALYSDRQSIL 435

10

20

Query: 411 QEANI--TIQLFINREFGSLGAINV 433

N+ +IQ+ I R G+ G + V

Sbjct: 436 IGQNLIRSIQINITRLAGTFGDVAV 460

15 Score = 230 (34.5 bits), Expect = 2.3e-14, Sum P(3) = 2.3e-14 Identities = 98/368 (26%), Positives = 164/368 (44%)

Query: 1240 EMEQD-

SLVTLNVERLKGTYGRITIAWEADGSISDIFPTSGVILFTEGQVLSTITLTILA 1298
E+E+D L+ + V RL GYYG +T + + S + P GV

ST+T

Sbjct: 71 EVEEDVGLIMIPVVRLHGTYGYVTADFISQSSSAS--P-GGVDYILHG---STVTFQH-G 123

25 Query: 1299 DNIPELSEVVIVTLTRITTEGVEDSYKGATIDQDRSKSVITTL--- PNDSPFGLVGWRAA 1355

N+ ++ +I E +E GAT + +++ +

6+

+DSPFG++ +

Sbjct: 124

30 QNLSFINISIIDDNESEFEEPIEILLTGATGGAVLGRHLVSRIIIAKSDSPFGVIRFLNQ 183

Query: 1356 SVFIRVAEPKENTTTLQLQIARDKGLLGDIAIHLRAQ-PNFLLHVDNQATENEDYVLQET 1414

S I +A P +T L L + R GLLG+I ++ PN + Q

35 D V

Sbjct: LB4 SK-ISIANPN-

PES 32--V9GAIGRN99LLABD2N9DVTBWNVBEIQEALLPQNRDIADPV--SG 239

Query: 1415 IIIMKENIKEAHAEV-

40 SILPDDLPELEEGFIVTITEVNLVNSDFSTGQPSVRRPGMEIAE 1473 + E + +I P + E+EE FI+ +++LV

++
Sbjct: 240 LFYFGEGEGGVRTIILTIYPHEEIEVEETFII---KLHLVK---GEAKLDSRAKDVT- 290

45

Query: 1474

+ 2 + WYV

50 Sbjct: 291

LTIQEFGDPNGVVQFAPETLSKKTYSEPLALEGPLLITFFVRRVKGTFGEIMVYWELSSE 350

Query: 1534

SAGLEDFKPSHGILEFADKQVTAMIEITIIDDAEFELTETFNISLISVAGGGRLGDDVVV 1593
55 EDF + G AD + A ++ ++ D E+ E + I L+SV GG

Sbjct: 351

FDITEDFLSTSGFFTIADGESEASFDVHLLPDEVPEIEEDYVIQLVSVEGGAELDLEKSI 410

Query: 1594 T-VVIPQNDSPFGVF 1607 + ND P GVF. T 411 TWFSVYANDDPHGVF 425 Sbict: 5 Score = 190 (28.5 bits), Expect = 7.5e-11, Sum P(3) = 7.5e-11 Identities = 136/591 (23%), Positives = 247/591 (41%) 67 SGRAVAGSDYEPVTRQWAIMREGDEFANLTVSILPDDFPEMDESFLISLLEVHLMNISAS 126 10 +G A D+EPV Q+ + ++I+ D E++E F I+L Sbjct: AGSAKPLEDFEPVQNGELFFQKFQTEVDFEITIINDQLSEIEEFFYINLTSVEIRGLQKF 766 15 127 LKN-QPTIGQP-NISTVVIALNGDAFGVFVIY-Query: SISPNTSEDGLFVEVQEQPQTLVELMI 183 NP+ +++ + I N D G+ + + + D 20 Sbict: 767 DVNWSPRLNLDFSVAVITILDNDDLAGMDISFPETTVAVAVDTTLIPVETESTTY--LST 824 184 HRTGGSLGQVAVEWRVVGGTATEGLDFIGAGEILTF--AEGETKKTVILTXXXXXXXXXX 241 25 L +V T G+ I ++K 825 SKTTTILQPTNVV-AIV--TEATGVSAIPE-KLVTLHGTPAVSEKPDVATVTANVSIHGT AAD 242 XXXXXXLVYTEGGSRILPSSDTVRVNILANDNVAGIVSF--30 QTASRSVIGHEGEILQFHV 299 +VY E + +T V I G VS LF Sbict: 881 FSLGPSIVYIEEEMKN-GTFNTAEVLIRRTGGFTGNVSITVKTFGERCAQMEPNALPF-- 937 35 Query: 300 IRTFPGRGNVTVNWKIIGQNLELNFANFSGQLFFPEGSLNTTLFVHLLDDNIPEEKEVYQ 359 R G N+T W + E +F + L F +G + V + LDD +PE +E + 40 Sbjct: 938 -RGIYGISNLT--WAVE----EEDFEEQTLTLIFLDGERERKVSVQILDDDEPEGQEFFY 990 360 VILYDVRTQGVPPAGIALLDAQ---GYAA--Querv: VLTVEASDEPHGVLNFALSSRFVL-LQEA 413 45 P G +++ + V L + G+AA ++ + ZDL L+E Sbjct: 991 VFLTN----PAGGARIVEGKDDTGFAAFAMVIITGSDLHNGIIGFSEESASGLELREG 1044 50 414 NITIQLFI-----NREFGSLGAI-NVTYTTVPGMLSLKNQTVGNLAEPEVDFVPIIGFL 466

E

+ +L + NR F + VT ++ L+ V NL E E+ Sbjct: 1045 AVMRRLHLIVTRQPNRAFEDVKVFWRVTLNKT--VVVLQKDGV-NLME-55 ELQSVS--GTT 1098

467 ILEEGETAAAINITILEDDVPELEEYFLVNL--TYVGLTMAASTSFPPRLDSEGLTAQVI 524

45

Query: 1289 LSTITLTILADNIPELS-EVVIVTLTRITTEGVEDSYK---GATIDQDRSKSVITTLPND 1344

50 +2 E+ +2 Т GA I+ I + DSbict: 1094 VSGTTTCTMGQTKCFISIELKPEKVPQVEVYFFVELYEATAGAAINNSARFAQIKILESD 1153

55 Query: 1345 SPFGLVGWRAASVFIRVAEPKENTTTLQLQIARDKG--LLGDIAI---HLRAQPNFLLHV 1399 Z LV + R+A T + LQ + ARD G L + +LR+

WO 01/98454 PCT/IB01/02050 Sbjct: 1154 ESQSLVYFSVGS---RLAVAHKKATLISLQVARDSGTGLMMSVNFSTQELRSAETIGRTI 1210 5 +D+V+ E ++ + + A + +V + P+ Sbict: 1211 ISPAISGKDFVITEGTLVFEPGQRSTVLDVILTPE 1245 Score = 186 (27.9 bits), Expect = 2.5e-13, Sum P(3) = 2.5e-13Identities = 75/242 (30%), Positives = 113/242 (46%) 10 Query: 1206 EIGINDSVTITILSNDDAYGIVAFARNSLYKRVEEMERDSLVTLNVERLKGTYGRITIAW 1265 VII+ ND+A GI+ F + Y E E L+ + V RL GTYG +T 40 EIGNISIVRIIIMKNDNAEGIIEF--15 Sbict: DPKYTAFEVEEDVGLIMIPVVRLHGTYGYVTADF 97 Query: 1266 EADGSIS----DIFPTSGVILFTEGQVLSTITLTILADNIPELSEVVIVTLTRITTEGV 1320 20 + 2 + D + F GQ LS I ++I+ DN LT T G Sbjct: 98 ISQSSSASPGGVDYILHGSTVTFQHGQNLSFINISIIDDNESEFEEPIEILLTGAT--G- 154 25 Query: 1321 EDSYKGATIDQDRSKSVITTLPNDSPFGLVGWRAASVFIRVAEPKENTTTLQLQIARDKG 1380 +DSPFG++ + S I +A P +T L GA + + +I L + R GSbjct: 155 -----GAVLGRHLVSRIIIA-KSDSPFGVIRFLNQSK-ISIANPN-30 STMILSLVLERTGG 206 Query: 1381 LLGDIAIHLRAQ-PNFLLHVDNQATENEDYVLQETIIIMKENIKEAHAEV-SEPT 347004112 LLG+I ++ PN 35 +IP+E207 LLGEIQVNWETVGPNSQEALLPQNRDIADPV--Sbjct: SGLFYFGEGEGGVRTIILTIYPHEEIE 264 1439 LEEGFIVTI 1447 Query: 40 +EE FI+ + Sbjct: 265 VEETFIIKL 273 Score = 179 (26.9 bits), Expect = 9.4e-34, Sum P(3) = 9.4e-34Identities = 65/244 (26%), Positives = 114/244 (46%) 45 Querv: 581 DYIFTPMILHFADGERYKNVNIMILDDDIPEGDEKFQLILTNPSPGLEL--GKN----T 633 + L F DGER + V++ ILDDD PEG E F + LTNP G ++ GK+ 50 Sbjct: DFEEQTLTLIFLDGERERKVSVQILDDDEPEGQEFFYVFLTNPQGGAQIVEGKDDTGFAA 1013 634 IALIIVLANDDGPGVLSFNNSEHFFLREPTALYVQESVAVLYIVREPAQG--Querv: ---LFGTV LAB

L + R+P +

G++ F+

1014 FAMVIITGSDLHNGIIGFSEESQSGLELREGAVMRR--

A++I+ +D

LHLIVTR@PNRAFEDVKVFWRV 1071

55

Sbict:

Query: LAS TVQ--FIVTEVNSSNESKDLTPSKGYIVLEEGVRFKALQISAILDTEPEMDEYFVCTLFN 746 +V + + N + +LG . G + I Sbict: 1072 TLNKTVVVLQKDGVNLMEELQSVSGTTTCTMGQTKCFISIELKPEKVPQVEVYFFVELYE 1131 747 PTGGARLGVHVQ-10 TLITVLQNQAPLGLFSISAVENRATSIDIEEANRTVYLNVSRTNGID 805 T GA + + I +L++ L S V +R ++ ++A V+R +6 . Sbjct: 1132 ATAGAAINNSARFAQIKILESDESQSLVYFS-VGSRL-AVAHKKAT-LISLQVARDSGTG 1188 15 BOL LAVSVQUET B14 Query: L + VZ + T1189 LMMSVNFST 1197 Sbjct: 20 Score = 174 (26.1 bits), Expect = 4.1e-32, Sum P(3) = 4.1e-32 Identities = 58/200 (29%), Positives = 102/200 (51%) Query: 1159 DFIPSSGELIFEPGEREATIAVNILDDTVPEKEESFKVQLKNPKGGAEIGINDSVT-ITI 1217 25 GE EA+ V++L D VPE EE + +QL + +GGAE+ + DF+ +SG ++ T+2 Sbjct: 356 DFLSTSGFFTIADGESEASFDVHLLPDEVPEIEEDYVIQLVSVEGGAELDLEKSITWFSV 415 30 Query: L218 LSNDDAYGIVAFAQNSLYKQVEEMEQDSL--VTLNVERLKGTYGRITIAWEADGSISDIF 1275 +NDD +G+ A +Q + Q+ + + +N+ RL GT+G + + 7.0 416 YANDDPHGVFALYSD---Sbict: 35 RQSILIGQNLIRSIQINITRLAGTFGDVAVGLRIS---SDHK 469 Querv: 1276 PTSGVILFTEGQVLSTITLTILADNIPELSEVVI----VTLTRITTEGVEDSYKGA-TI 1329 E Q++ T D +P ++V + TL +T 40 + G TI Sbjct: EQPIVTENAERQLVVKDGATYKVDVVPIKNQVFLSLGSNFTLQLVTVMLVGGRFYGMPTI 529 Query: 1330 DQDRSKSVITTLPNDSPFGLVGWRAAS 1356 45 Q+ +KS + + + VG+ + + 530 LQE-AKSAVLPVSEKAANSQVGFESTA 555 Sbjct: Score = 145 (21.8 bits), Expect = 4.3e-24, Sum P(3) = 4.3e-24 Identities = 104/396 (26%), Positives = 170/396 (42%) 50 Query: 88 EGDEFANLTVSILPDDFPEMDESFLISLLEVHLMNISASLKNQPTIGQPNISTVVIALNG 147 +G+ A+ V +LPD+ PE++E ++I L+ V A L + +I

368 DGESEASFDVHLLPDEVPEIEEDVYVQLVSVEG---GAELDLEKSI-----

55

Sbjct:

TUFSVYAND 419

Query: 148

DAFGVFVIYSISPNTSEDGLFVEVQEQPQTLVELMIHRTGGSLGQVAVEWRVVGGTATEG 207 D GVF +YS +++ I R G+ G VAV R+ . **D** + + +

5 420 DPHGVFALYS----Sbjct: DRQSILIGQNLIRSIQINITRLAGTFGDVAVGLRISSDHKEQP 472

208 LDFIGAGEILTFAEGETKKTVILTXXXXXXXXXXXXXXXXXLVYTE-GGSRI-E45 TG-ZZ91-

10 +G T K ++ LV : GR

Sbjct: 473

IVTENAERQLVVKDGATYKVDVVPIKNQVFLSLGSNFTLQLVTVMLVGGRFYGMPTILQE 532

15 264 VRVNIL-ANDNVAGI-VSFQTASRSVIGHEGEILQFHVIRTFPGR-Query: GNVTVNWKI-IGQN 319

> V F++ + ++ + +L. ++ A HV++GG++V

Ы Sbict: -- ZTDATINMIDFATZEFDVDZNAAMEZVGLVAZNA 533 AKSAVLPVSEKAANSQVGFESTAFQLMNITAGTS--

20 HVMISRRGTYGALSVAUTTGYAPG 590 <

> 350 FEF----Query: NFANFSGQLFFPEGSLNTTLFVHLLDDNIPEEKEVYQVILYDVRTQGVPP 372 LE+ N GLFGР E + + L

25 Sbjct: 591 LEIPEFIVVGNMTPTLGSLSFSHGEQRKGVFLWTFPS--PGWPEAFVLHLSGVQSSA--P 646

Query: 373 ·

30 AGIALLDAQGYAAVLTVEASDEPHGVLNFALSSRFVLLQEANITIQLFINREFG-SLGAI 431 G L G+ + A EP GV F+ SSR +++ E FG 647 GGAQL--RSGF----

IVAEIEPMGVFQFSTSSRNIIVSEDTQMIRLHVQRLFGFHSDLI 699

35

432 NVTYTTVPGMLS-LKN-QTV--GNLA----EPEVDF-Query: VPIIGFLILEEGETAAAINITIL 482

V+Y T G L++ + V G L + EVDF + II LEE

IN+T+

40 Sbjct: 700 KVSYQTTAGSAKPLEDFEPVQNGELFFQKFQTEVDFEITIINDQ-LSEIEEFFYINLTSV 758

Query: 483 E 483

45 Sbjct: 759 E 759

> Score = 142 (21.3 bits), Expect = 5.6e-05, Sum P(3) = 5.6e-05Identities = 54/175 (30%), Positives = 76/175 (43%)

50 Query: 1435

DLPELEEGFIVTITEVNLVNSDFSTGQPSVRRPGMEIAEIMIEENDDPRGIFMFHVTRGA 1494 G+ TIEN + D QΡ + I I+I +ND+

Sbict: 16 DLYDFGRGYDFTIQE-NGLQID----QPP-

55 EIGNISIVRIIIMKNDNAEGIIEFDPK--- LL

> Querv: 1495 GEVITAYEXXXXXXXXXXXXXXAGSFGAVNVYW--KASPDSAGLEDFKPSHGILEFADK 1552

WO 01/98454 PCT/IB01/02050 TA+E G++G V ++Z Z G D+ + F Sbjct: 67 ---YTAFEVEEDVGLIMIPVVRLHGTYGYVTADFISQSSSAGPGVDYILHGTVTFQHG 123 5 Query: 1553 QVTAMIEITIIDDAEFELTETFNISLISVAGGGRLGDDVVVTVVIPQNDSPFGVFGF 1609 Q + I I+IIDD E E E I L GG LG +V ++I ++DSPFGV F 124 10 Sbjct: QNLSFINISIIDDNESEFEEPIEILLTGATGGAVLGRHLVSRIIIAKSDSPFGVIRF 180 Score = 125 (18.8 bits), Expect = 4.0e-25, Sum P(3) = 4.0e-25Identities = 77/308 (25%), Positives = 134/308 (43%) 15 Querv: 1141 LVGAHSHIYELAYISSHS-----DFIP-SSGELIFEPGEREATIAVNILDDTVPEKEES 1193 DF P +GEL F+ + E L G HS + +++Y ++ + I++D + E EE 20 Sbict: 691 LFGFHSDLIKVSYQTTAGSAKPLEDFEPVQNGELFFQKFQTEVDFEITIINDQLSEIEEF 750 Query: 1194 FKVQLKNP--KGGAEIGINDSVTITILSNDDAYGIVAFAQNSLYKQVEEMEQDSLVTLNV 1251 25 F + L + +G + +NS + +D + ++Sbjct: 751 FYINLTSVEIRGLQKFDVNWSPRLNL---DFSVAVITILDN----DDLAGMDI 796 Query: 1252 30 ERLKGTYGRITIAWEADGSISDIFPTSGVILFTEGQVLSTITLTILADNIPELSEVVIVT 1311 + T + + + T + Z + + + C A + T +++ + F 797 ----SFPETTVAVAVDTTLIPVETESTTYLSTS-Sbjct: KTTTILQPTNVVAIVTEATGVSAIP 850

35

Query: 1312 LTRITTEGVEDSYKGATIDQDRSKSVITTLPNDSPFGLVGWRAASVFIRVAEPKENT-TT 1370 D Z N T V +T G Т + V+I

KTT 40 Sbjct: 851 EKLVTLHG-----TPAVSEKPDVATVTANVSIHGTFSLGPSIVYIE-EEMKNGTFNT 901

Query: 1371 LQLQIARDKGLLGDIAIHLRA-----QPNFL----LHVDNQ--ATENEDYVLQETI 1415

++ I R G G+++I ++ +PN L + N A E ED+ Q

Sbjct: 902

50

AEVLIRRTGGFTGNVSITVKTFGERCAQMEPNALPFRGIYGISNLTWAVEEEDFEEQTLT 961

J4J6 IIMKENIKEAHAEVSILPDDLPELEEGFIVTIT J448 Query: V IL DD PE +E F V +T +I + +E Sbict: 962 LIFLDGERERKVSVQILDDDEPEGQEFFYVFLT 994

Score = 123 (18.5 bits), Expect = 6.0e-28, Sum P(3) = 6.0e-2855 Identities = 91/372 (24%), Positives = 150/372 (40%)

		÷	

WO 01/98454 PCT/IB01/02050 386 VLTVEASDEPHGVLNFALSSRFVLLQEA--NITI---Query: QLFINREFGSLGAINVTYTTV-- 438 V TV A+ F+L V ++E NT ++ I R +++T T 868 VATVTANVSIHGT --Sbict: FSLGPSIVYIEEEMKNGTFNTAEVLIRRTGGFTGNVSITVKTFGE 925 439 -----PGMLSLKN-QTVGNL--AEPEVDFVPIIGFLILEEGETAAAINITILEDDVPEL 489 10 PL+ + NL A E DF LI +GE IL+DD PE Sbjct: 926 RCAQMEPNALPFRGIYGISNLTWAVEEEDFEEQTLTLIFLDGERERKVSVQILDDDEPEG 985 490 EEYFLVNLTYVGLTMAASTSFPPRLDSEGLTA--QVIIDANDGARGVI---15 Query: EWQQSRFEV 544 +E+F V LT D GA VII +D G+I E QS E+ Sbjct: 986 QEFFYVFLT----20 NP@GGA@IVEGKDDTGFAAFAMVIITGSDLHNGIIGFSEES@SGLEL 1041 545 NE--THGSLTLVAQRS-REPLGHVSLF--Query: VYAQNLEAQVGLDYIFTPMILHFADGERYKN 599 L L+ R E 25 Sbjct: 3042 REGAVMRRLHLIVTRQPNRAFEDVKVFWRVTLNKTVVVLQKDGVNLMEELQSVSGTTTCT 1101 LOO -----VNIMILDDDIPEGDEKFQLILTNPSPGLELGKNT-30 IALIIVLANDDGPGVLSF 651 ++I + + + P + + F + L+ 6 A I +L +D+ ++ F Sbict: 1102 MGQTKCFISIELKPEKVPQVEVYFFVELYEATAGAAINNSARFAQIKILESDESQSLVYF lll. 35 L52 NNSEHFFLREPTALYVQESVAVLYIVREPAQGLFGTVTVQFIVTEVNSSNE--SKDLTPS 709 L + R+ GL ++V F E+ S+ ++P+ 40 Sbjct: 1162 SVGSRLAVAHKKATLIS----LQVARDSGTGLM--MSVNFST@ELRSAETIGRTIISPA 1214 710 ---KGYIVLEEGVRFKALQISAILD 731 Query: K +++ E + F+ Q S +LD45 Sbict: 1215 ISGKDFVITEGTLVFEPGGRSTVLD 1239 Score = 120 (18.0 bits), Expect = 1.8e-22, Sum P(3) = 1.8e-22Identities = 77/316 (24%), Positives = 127/316 (40%) 50 Query: 1255 KGTYGRITIAWE---ADGS-----ISDIFPTSGVILFTEGQVLSTITLTILADNIPEL 1304 + ++ PT G + F+ G+ +GTYG +++AW A G Sbjct: 573 55 RGTYGALSVAUTTGYAPGLEIPEFIVVGNMTPTLGSLSFSHGEQRKGVFLUTFPS--PGW 630 Query: 1305

```
GV+S G
                  E ++ L+
                                          Q RS ++
                                                         P G+ +
    I V+E
    Sbjct:
           L31 PEAFVLHLS----GVQSSAPGGA--QLRSGFIVAEI---
    EPMGVFQFSTSSRNIIVSE- 679
 5
    Query: 1365 KENTTTLQLQIARDKGLLGDIAIHLRAQPNFLLHVDNQATENEDYV-
    LQETIIIMKENIK 1423
                   +T ++L +R G
                                   D+I+Q
                                                          ED++Q
    Sbjct:
10
             68D --DTQMIRLHVQRLFGFHSDL-IKVSYQTTA----
    GSAKPLEDFEPVQNGELFFQKFQT 731
    Query: 1424 EAHAEVSILPDDLPELEEGFIVTITEVNLVN-
    SDFSTGQPSVRRPGMEIAEIMIEENDDP 1482
15
                    E++I+ D L E+EE F + +T V +
    I +NDD
    Sbjct:
             732
    EVDFEITIIND@LSEIEEFFYINLTSVEIRGL@KFDVNWSPRLNLDFSVAVITILDNDDL 791
20
    Query: 1483 RGI-FMFHVTRGAGEVITAY---
    EXXXXXXXXXXXXXAGSFGAVNVYWKASPDSAGLE 1538
                  G+
                     FTAVT
    AZ +A+
    Sbjct:
             792
25
    AGMDISFPETTVAVAVDTTLIPVETESTTYLSTSKTTTILQPTNVVAIVTEATGVSAIPE A51
    Query: 1539 DFKPSHGILEFADKQVTAMIEITIIDDAEFEL 1570
                     HG
                           ++K A + + '
    Sbjct: 852 KLVTLHGTPAVSEKPDVATVTANVSIHGTFSL 883
30
     Score = 113 (17.0 bits), Expect = 9.4e-34, Sum P(3) = 9.4e-34
     Identities = 28/87 (32%), Positives = 50/87 (57%)
    Query: 1156 SHSDFIPSSGELIFEPGEREATIAVNILDDT--
35
    VPEKEESFKVQLKNPKGGAEIG-INDS 1212
              S DF+ + G L+FEPG+R
                                      + V + +T +
                                                      + F++ L +PKGGA
    Sbjct: 1216
    SGKDFVITEGTLVFEPGQRSTVLDVILTPETGSLNSFPKRFQIVLFDPKGGARIDKVYGT 1275
40
            1573 ALLILIAND TAKEN STAND TAKEN 1540
                  + V + + J + A + L + + V +
            J276 ANITLVSDADSQAIWGLA-DQLHQPVND J302
    Sbjct:
     Score = 93 (14.0 bits), Expect = 4.1e-32, Sum P(3) = 4.1e-32
45
     Identities = 57/222 (25%), Positives = 90/222 (40%)
    Query: 1404 TENEDYVL--QETIIIMKENIKEAHAE---VSILPDDLPEL------
    EEGFIVTITEVN 1451
50
                TE+
                          + T I+
                                   N+
                                          Ε
                                              VS +P+ L L
                                                                E+
    + T+T
    Sbict:
    TESTTYLSTSKTTTILQPTNVVAIVTEATGVSAIPEKLVTLHGTPAVSEKPDVATVTANV A75
55
    Query: 1452 LVNSDFSTGQPSVRRPGMEIAEIMIEENDDPRGIFMFHVTRGAGEV-
    ITAYEXXXXXXXX 1510
                                 + I E M
                 ++ FS G PS+
                                                           G V IT
```

Na.			
4.5			
÷			

Sbjct: 876 SIHGTFSLG-PSI---VYIEEEMKNGTFNTAEVLIRRTGGFTGNVSITVKTFGERCAQM 930

Query: 1511

5 XXXXXXXAGSFGAVNVYWKASPDSAGLEDFKPSHGILEFADKQVTAMIEITIIDDAEFEL 1570
G +G N+W EDF+ L F D + + +

I+DD E E

Sbjct:

Sbjct: 931 EPNALPFRGIYGISNLTWAVEE----EDFEEQTLTLIFLDGERERKVSVQILDDDEPEG 985

10

Query: 1571 TETFNISLISVAGGGRL--GDD-----VVVTVVIPQNDSPFGVFGFEEKTVS 1615

E F + L + GG ++ G D V+I +D G+ GF E++ S
986 QEFFYVFLTNPQGGAQIVEGKDDTGFAAFAMVIITGSDLHNGIIGFSEESQS

15 1037

Score = 93 (14.0 bits), Expect = 1.0e-18, Sum P(3) = 1.0e-18 Identities = 51/238 (21%), Positives = 107/238 (44%)

20 Query: 600 VNIMILDDDIPEGDEKFQLILTNPSPGLELGKNTIALIIVLANDDGPGVLSFNNSEHFF 658
++I + ++P+ + F + L + G + + A I +L +D+ ++
F+

Sbjct: 1109

25 ISIELKPEKVPQVEVYFFVELYEATAGAAINNSARFAQIKILESDESQSLVYFSVGSRLA 1168

Query: L59 LREPTALYVQESVAVLYIVREPAQGLFGTVTVQFIVTEVNSSNE-SKDLTPS---KGYI 713
+ A + L + R+ GL ++V F E+ S+

30 +++ K ++
Sbjct: 11b9 VAHKKATLIS----LQVARDSGTGLM-MSVNFSTQELRSAETIGRTIISPAISGKDFV 1221

Query: 714 VLEEGVRFKALQISAILDT--EPE---MDEY---FVCTLFNPTGGARLG-

35 VHVQTLITYL 764

+ E + F+ Q S +LD PE ++ + F LF+P GGAR+ V+

Sbjct: 1222

ITEGTLVFEPGQRSTVLDVILTPETGSLNSFPKRFQIVLFDPKGGARIDKVYGTANITLV 1281

40

Query: 765 QNQAPLGLFSISAVENRATSIDI-EEANRTVYLNVSRTNGIDLAVSVQWETVSETAFGMR 823 + ++ ++ ++ DI T+ +V+ T D +S

45 Sbjct: 1282 SDADSQAIUGLADQLHQPVNDDILNRVLHTISMKVA-TENTDEQLSAMMHLIEKIT--TE 1338

Query: 824 GMDVVFSV 831

G FSV

50 Sbjct: L339 GKIQAFSV L346

Score = 92 (13.8 bits), Expect = 9.5e-25, Sum P(3) = 9.5e-25 Identities = 44/177 (24%), Positives = 82/177 (46%)

55 Query: 680
PAQGLFGTVTVQFIVTEVNSSNESKDLTPSKGYIVLEEGVRFKALQISAILDTEPEMDEY 739
P+G++G + + V E + E + LT ++ +G R + + + + I
EPE E+

PCT/IB01/02050 Sbict: 936 PFRGIYGISNLTWAVEEEDF--EEQTLT----LIFLDGERERKVSVQILDDDEPEGQEF 988 Query: 740 FVCTLFNPTGGARL-----GVHVQTLITVLQNQAPLGLFSISAVENRATSIDIEEAN- 791 F L NP GGA++ Ε Sbjct: 989 FYVFLTNPQGGAQIVEGKDDTGFAAFAMVIITGSDLHNGIIGFS--EESQSGLELREGAV 1046 10 792 -RTVYLNVSRT-NGIDLAVSVQWE-TVSETAF----Query: GMRGMDVVFSVFQSFLDESASGW 843 R ++L V+R V V W T+++T V2 + +M Sbjct: 1047 MRRLHLIVTRQPNRAFEDVKVFWRVTLNKTVVVLQKDGVNLMEELQSVSGTTTCTMGQTK 1106 15 Query: 844 CFFTLE 849 CF ++E Sbjct: 1705 CLIZIE 7775 20 Score = 91 (13.7 bits), Expect = 6.6e-32, Sum P(3) = 6.6e-32 Identities = 49/153 (32%), Positives = 70/153 (45%) Query: 1466 RPGMEIAEIMIEENDDPRGIFMFHVTRGAGEVITAYEXXXXXXXXXXXXXXXXXAGSFGAVN 1525 25 R G +AEI +P G+F F + + +I + + Sbjct: L52 RSGFIVAEI----EPMGVFQFSTS--SRNIIVSEDT@MIRLHV@RLFGFHSD---LIK 700 30 Query: 1526 VYWKASPDSAG-LEDFKP-SHGILEFADKQVTAMIEITIIDDAEFELTETFNISLISVAG 1583 V ++ + SA LEDF+P +G L F EITII+D E+ E F I+L SV 35 Sbjct: 701 VSYQTTAGSAKPLEDFEPVQNGELFFQKFQTEVDFEITIINDQLSEIEEFFYINLTSVEI 760 1584 GG-----RLGDDVVVTVV-IPQNDSPFGV-FGFEEKTVS 1615 RL D V V+ I ND G+ 40 Sbjct: 761 RGLQKFDVNWSPRLNLDFSVAVITILDNDDLAGMDISFPETTVA 804 Score = 65 (9.8 bits), Expect = 8.8e-29, Sum P(3) = 8.8e-29 Identities = 26/99 (26%), Positives = 50/99 (50%) 45 - duery: 1232 NSLYKQVEEMEQDSLVTLNVERLKGTYGRITIAWEADGS----ISDIF--PTSGVILFTE 1285 K+ + + ++++ GT $\cdot IT + + AD$ ++D+ IL 1250 NSFPKRFQIVLFDPKGGARIDKVYGT-50 ANITLVSDADSQAIWGLADQLHQPVNDDIL--- 1305 Query: 1286 GQVLSTITLTILADNIPELSEVVIVTLTRITTEGVEDSYKGAT 1328 +VL TI++ + +N E ++ + + TTEG 1306 NRVLHTISMKVATENTDEQLSAMMHLIEKITTEGKIQAFSVAS 1348 Sb ict: 55

WO 01/98454

Score = 48 (7.2 bits), Expect = 1.9e-27, Sum P(3) = 1.9e-27

Identities = 23/115 (20%), Positives = 44/115 (38%)

PRD



Query: 1499 TAYEXXXXXXXXXXXXXXXAGSFGAVNVYWKAS-----PDSAGLEDFKPSHGILEFAD 1551 G++GA++V W P+ TA++ G L F+ 5 Sbjct: 554 TAFQLMNITAGTSHVMISRRGTYGALSVAUTTGYAPGLEIPEFIVVGNMTPTLGSLSFSH 613 Query: 1552 KQVTAMIEITIIDDAEFELTETFNISLI--SVAGGGRLGDDVVVTVVIPQNDSPFGVFGF 1609 + 2++ 2 66 +1 10 P GVF F 614 GERRKGVFLWTFPSPGWPEAFVLHLSGV&SSAPGGARLRSGFIVAEI----Sbict: EPMGVFQF 668 15 Pedant information for DKFZphamy2_10p7, frame 3 ______ Report for DKFZphamy2_10p7.3 20 ELENGTHI 1615 EMWI 177600-58 [DI] 4.37 TREMBL: AF055084_1 gene: "VLGR1"; product: "very 25 ELOMOLE large G-protein coupled receptor-1"; Homo sapiens very large Gprotein coupled receptor-1 (VLGR1) mRNA, complete cds. 5e-24 EBFOCKZI BLOTAL3V EBLOCKSB BLOO713B Sodium:dicarboxylate symporter family proteins 30 EBF0CK21 AEOOLOSS [BLOCKZ] PR00412C EBFOCKZ3 BF00954E **EPIRKWD** heart le-D8 ion transport le-08 35 **EPIRKU** transmembrane protein 3e-08 **EPIRKWI** phosphoprotein 2e-08 CPIRKW1 membrane protein le-08 **EPIRKUI** [PROSITE] MULTICOPPER_OXIDASEL All_Beta 40 EKW1 LOW_COMPLEXITY 2.60 % EKWI SEQ DAWADAWALYTCATLCLKEQACSAFSFFSASEGPQCFWMTSWISPAVNNSDFWTYRKNMTxxxxxxxxxx...... 45 SEG ccchhhhhhhhhhhhhhhhhhhheeeeeccccceeeeecccee PRD RVASLFSGQAVAGSDYEPVTRQWAIMQEGDEFANLTVSILPDDFPEMDESFLISLLEVHL SEQ SEG 50 PRD MNISASLKNQPTIGQPNISTVVIALNGDAFGVFVIYSISPNTSEDGLFVEVQEQPQTLVE ZEQ SEG hccccccccccccceeeeeecccceeeeeeccccceee PRD 55 LMIHRTGGSLGQVAVEWRVVGGTATEGLDFIGAGEILTFAEGETKKTVILTILDDSEPED SEQ SEG

	SEQ	DESIIVSLVYTEGGSRILPSSDTVRVNILANDNVAGIVSFQTASRSVIGHEGEILQFHVI
_	PRD	ccceeeeeeccccccccceeeeeecccceeeeeeecccceeee
5	SEQ	RTFPGRGNVTVNWKIIGQNLELNFANFSGQLFFPEGSLNTTLFVHLLDDNIPEEKEVYQV
•	SEG	•••••
	PRD	ecccccceeeeeeeeccccccccceee
10	SEQ	ILYDVRTQGVPPAGIALLDAQGYAAVLTVEASDEPHGVLNFALSSRFVLLQEANITIQLF
	SEG PRD	eeccceeeccchhhhhhhhccccceeeeeecccccceeeeecececccceeee
	LIA	secressecrumnuluuccccsssssssssssssssssssssssssssssss
	SEQ	INREFGSLGAINVTYTTVPGMLSLKNQTVGNLAEPEVDFVPIIGFLILEEGETAAAINIT
15	SEG .	
	רוע	CCCCCCCeeeeeeecccccccccccccceee
	SEQ	ILEDDVPELEEYFLVNLTYVGLTMAASTSFPPRLDSEGLTAQVIIDANDGARGVIEWQQS
20	SEG PRD	ecccchhhhhheeeeeeecceeeccccccccccceeeeee
20	rkv.	eccccunuuuuneeeeeeecccecccccccccccccccc
	SEQ	RFEVNETHGSLTLVAQRSREPLGHVSLFVYAQNLEAQVGLDYIFTPMILHFADGERYKNV
	SEG PRD	eeeeccccceeeeecccccceeeeeecccccccccccc
25	LIV	
	SEQ	NIMILDDDIPEGDEKFQLILTNPSPGLELGKNTIALIIVLANDDGPGVLSFNNSEHFFLR
	SEG PRD	PREPRICTOR CONTRACTOR
•	FILD	54666CCCCCCCCG666666CCCCCCCCCG666666CCCG66666
30	SEQ	EPTALYV@ESVAVLYIVREPA@GLFGTVTV@FIVTEVNSSNESKDLTPSKGYIVLEEGVR
	SEG PRD	ccceeeccchhhhhhhhccccceeeeeeeeeecccccccc
	IND	
25	SEQ	FKALQISAILDTEPEMDEYFVCTLFNPTGGARLGVHVQTLITVLQNQAPLGLFSISAVEN
35	SEG PRD	eeeeeeecccchhhhhhheeeeccccceeehhhhhhhhh
	SEG	RATSIDIEEANRTVYLNVSRTNGIDLAVSVQWETVSETAFGMRGMDVVFSVFQSFLDESA
40	PRD	hhhhhcccccceeeeeecccchhhhheeeeeccceeeecccceeeeeccccc
	SEQ	SGWCFFTLENLIYGIMLRKSSVTVYRWQGIFIPVEDLNIENPKTCEAFNIGFSPYFVITH
	SEG PRD	CCeeeeccccceeecccceeecccceeeeccccceeeccccc
45	1 112	
	SEQ	EERNEEKPSLNSVFTFTSGFKLFLVQTIIILESSQVRYFTSDSQDYLIIASQRDDSELTQ
	SEG PRD	hhhh
	FKV	hhhhhcccceeeeeecccceeeeccccceeeeccccceee
50	SEQ	VFRUNGGSFVLHQKLPVRGVLTVALFNKGGSVFLAISQANARLNSLLFRUSGSGFINFQE
	SEG	**************************************
	PRD	eeeeccceeeeeccccceeeeeeeccccceeeeehhhhhh
	SEQ	VPVSGTTEVEALSSANDIYLIFAKNVFLGDQNSIDIFIWEMGQSSFRYFQSVDFAAVNRI
55	SEG	
	PRD	eeccccceeeecccceeeeeeecccceeeeecccceeeee
	SEQ	HSFTPASGIAHILLIGQDMSALYCUNSERNQFSFVLEVPSAYDVASVTVKSLNSSKNLIA

WO 01/98454							PCT/IB01/02	2050
	SEG PRD	eeccccceeeeeee	cccceee		ccceeeeee	•		
5	SEQ SEG PRD	EECCCEEEEEEEEE			• • • • • • • • •	<i>.</i>	. .	
10	SEQ SEG PRD	CCCCCEEECCCCEEE						
15	SEQ SEG PRD	eeeeeeeccceeeee						
13	SEQ SEG PRD	CCEEEEEEEECCCCE						• • • • • • •
20	SEQ SEG PRD	LLGDIAIHLRAQPNFI						
25	SEQ SEG PRD	CC6666666666CCC	• • • • • • •					
30	SEQ SEG PRD	YEVPPPLNVLQVPVVIxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx	xx					• • • • • • • • •
35	SEQ SEG PRD	TIIDDAEFELTETFN					• • • • • • •	• • • •
	•					•	•	
		·	Prosite	for	DKFZphamy	2_10p?-3	3	
40	PSDO	D79 151->1·	72 MUL	TICOF	PER_OXIDA	ZEI	РТОСПОП	76

(No Pfam data available for DKFZphamy2_10p7.3)

WO 01/98454 DKFZphamy2_11d2

5 group: transmembrane protein

DKFZphamy2_lld2 encodes a novel 552 amino acid protein without similarity to known proteins.

- The novel protein contains 2 transmembrane regions.
 No informative Blast results; no predictive prosite; pfam or scope motife.
- The new protein can find application in studying the expression profile of amygdala-specific genes and as a new marker for amygdala cells.

unknown protein

20

Pedant: TRANSMEMBRANE 2

Sequenced by EMBL

25 Locus: /map="lbpl3.3"

Insert length: 2939 bp

Poly A stretch at pos. 2920, polyadenylation signal at pos. 2869

30 1 GGCGGGTGAG AGGCCGCGGC GGCAGGTCCA CCTGGGCTTG CGAAGGCACA 51 GATTCCCCGT CCACAGCTCA CGACCAGATG CACCAGCAGG AGTCCACATC 101 GAGGACGTCC TCCGGGCACT CCCACGACCA GTGACCAGGA GTTAAACTTT 151 GGGATGTGCC CGTGATGTTG GACCACAAGG ACTTAGAGGC CGAAATCCAC 2D1 CCCTTGAAAA ATGAAGAAAG AAAATCGCAG GAAAATCTGG GAAATCCATC 35 251 AAAAAATGAG GATAACGTGA AAAGCGCGCC TCCACAGTCC CGGCTCTCCC 301 GGTGCCGAGC GGCGGCGTTT TTTCTTTCAT TGTTTCTCTG CCTTTTTGTG 351 GTGTTCGTCG TCTCATTCGT CATCCCGTGT CCAGACCGGC CGGCGTCACA 401 GCGAATGTGG AGGATAGACT ACAGTGCCGC TGTTATCTAT GACTTTCTGG 45% CTGTGGATGA TATAAACGGG GACAGGATCC AAGATGTTCT TTTTCTTTAT 40 5D1 AAAAACACCA ACAGCAGCAA CAATTTCAGC CGATCCTGTG TGGACGAAGG 551 CTTTTCCTCT CCCTGCACCT TTGCAGCTGC TGTGTCGGGG GCCAACGGCA LD1 GCACGCTCTG GGAGAGACCT GTGGCCCAAG ACGTGGCCCT CGTGGAGTGT **L51 GCTGTGCCCC AGCCAAGAGG CAGTGAGGCA CCTTCTGCCT GCATCCTGGT** 7D1 GGGCAGACCC AGTTCTTTCA TTGCAGTCAA CTTGTTCACA GGGGAAACCC 45 751 TGTGGAACCA CAGCAGCAGC TTCAGCGGGA ATGCGTCCAT CCTGAGCCCT **BD1 CTGCTGCAGG TGCCTGATGT GGACGGCGAT GGGGCCCCAG ACCTGCTGGT** 851 TCTCACCCAG GAGCGGGAGG AGGTTAGTGG CCACCTCTAC TCCGGCAGCA 9D1 CCGGGCACCA GATTGGCCTC AGAGGCAGCC TTGGTGTGGA CGGGGAAAGT 951 GGCTTCCTCC TTCACGTCAC CAGGACAGGT GCCCACTACA TCCTCTTTCC
1001 CTGCGCAAGC TCCCTCTGCG GCTGCTCTGT GAAGGGTCTC TACGAGAAGG
1051 TGACCGGGAG CGGCGGCCCG TTCAAGAGTG ACCCGCACTG GGAGAGCATG
1101 CTCAATGCCA CCACCCGCAG GATGCTTTCC CACAGCTCTG GAGCAGTGCG
1151 CTACCTGATG CATGTCCCAG GGAACGCCGG TGCAGATGTG CTTCTTGTGG
1201 GCTCAGAGGC CTTCGTGCTG CTGGACGGGC AGGAGCTGAC GCCTCGCTGG 50 55

PCT/1B01/02050 WO 01/98454 10 15 20 25 30 **BLAST Results** 35 No BLAST result 40 Medline entries No Medline entry 45

Peptide information for frame 2

50 ORF from 2555 bp to 2839 bp; peptide length: 95 Category: questionable ORF Classification: unclassified

55

I MCCEYPKELA VECVFGSVCA LSVDTGAALS LKRPRAPGMA SACLSPSGAH
51 TPTPCHPMQD SPLCLAAPEA QGQPWCSVLL GPPRSQSLSF VAKAA

BLASTP hits

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_11d2, frame 2 5

TREMBL:MMIGCF_2 Mouse ig gamma2a-b(c57bl/b allele) c gene and secreted

tail., N = 1, Score = 73, P = 0.1

10

>TREMBL:MMIGCF_2 Mouse ig gamma2a-b(c57bl/b allele) c gene and secreted

tail.

15

Length = 334

HSPs:

Score = 73 (11.0 bits), Expect = 1.1e-01, P = 1.0e-01 Identities = 16/49 (32%), Positives = 27/49 (55%) 20

44 LSPSGAHTPTPCHPMQDSPLCLAAPEAQGQPWCSVLLGPPRSQSLSFVA 92 Query: + P T PC P+++ P C AAP+ G P SV + PP+ + ++

96 IEPRVPITQNPCPPLKECPPC-AAPDLLGGP--SVFIFPPKIKDVLMIS Sbjct:

25 141

Peptide information for frame 3

30

ORF from 165 bp to 1820 bp; peptide length: 552 Category: putative protein Classification: Transmembrane proteins unclassified

35

1 MLDHKDLEAE IHPLKNEERK SQENLGNPSK NEDNVKSAPP QSRLSRCRAA 51 AFFLSLFLCL FVVFVVSFVI PCPDRPASQR MWRIDYSAAV IYDFLAVDDI LOL NGDRIQDVLF LYKNTNSSNN FSRSCVDEGF SSPCTFAAAV SGANGSTLWE 151 RPVAQDVALV ECAVPQPRGS EAPSACILVG RPSSFIAVNL FTGETLWNHS 201 SSFSGNASIL SPLLQVPDVD GDGAPDLLVL TQEREEVSGH LYSGSTGHQI 251 GLRGSLGVDG ESGFLLHVTR TGAHYILFPC ASSLCGCSVK GLYEKVTGSG 301 GPFKSDPHWE SMLNATTRRM LSHSSGAVRY LMHVPGNAGA DVLLVGSEAF 40 351 VLLDGQELTP RWTPKAAHVL RKPIFGRYKP DTLAVAVENG TGTDRQILFL 401 DLGTGAVLCS LALPSLPGGP LSASLPTADH RSAFFFWGLH ELGSTSETET 451 GEARHSLYMF HPTLPRVLLE LANVSTHIVA FDAVLFEPSR HAAYILLTGP 45 501 ADSEAPGLVS VIKHKVRDLV PSSRVVRLGE GGPDSDQAIR DRFSRLRYQS 551 EA

50

BLASTP hits

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_11d2, frame 3 55

No Alert BLASTP hits found

SEG

PRD

MEM

ZEQ

55

Pedant information for DKFZphamy2_lld2, frame 2

Report for DKFZphamy2_lld2.2 5 **ELENGTHD** 95 9757-38 EWWI [[q] 64.4 **EBFOCK21** 10 PR00521E EKUI Alpha_Beta MCCEYPKELAVECVFGSVCALSVDTGAALSLKRPRAPGMASACLSPSGAHTPTPCHPMQD SEQ 15 PRD SPLCLAAPEAQGQPWCSVLLGPPRSQSLSFVAKAA PRD cccccccccceeeecccccchhhhhhccc 20 (No Prosite data available for DKFZphamy2_11d2.2) (No Pfam data available for DKFZphamy2_11d2.2) 25 Pedant information for DKFZphamy2_lld2, frame 3 Report for DKFZphamy2_11d2.3 30 552 **ELENGTHI** EMWI 59659.68 5.84 [pI] **EBLOCK2** 35 PR002116 BLOOZABC Tissue inhibitors of metalloproteinases **EBFOCK2** proteins PR00436A **EBFOCKZ** TRANSMEMBRANE [KW] LOW_COMPLEXITY 40 **EKM3** 8-15 % MLDHKDLEAEIHPLKNEERKSQENLGNPSKNEDNVKSAPPQSRLSRCRAAAFFLSLFLCL SEQ SEG 45 PRD MEM SEQ **FVVFVVSFVIPCPDRPASQRMWRIDYSAAVIYDFLAVDDINGDRIQDVLFLYKNTNSSNN** SEG 50 PRD MEM FSRSCVDEGFSSPCTFAAAVSGANGSTLWERPVAQDVALVECAVPQPRGSEAPSACILVG SEQ

RPSSFIAVNLFTGETLWNHSSSFSGNASILSPLLQVPDVDGDGAPDLLVLTQEREEVSGH

	wo	01/98454	PCT/IB01/02050
	SEG PRD MEM	ccceeeeeccccccccccccceeeecceeeccccccc	ccchhhhhhhhhhhcc
5	SEQ SEG PRD MEM	cccccccccccccccceeeeeecccccc	ccccceeeecccccc
10	SEQ SEG PRD MEM	ccccccccchhhhhhhhhcccccceeeccccccceee	eccceeeecccccc
15	SEQ SEG PRD MEM	ccchhhhhhcccccccccceeeeeecccccceeeeeeccc	XXXXXXXXXXX
20	SEQ SEG PRD MEM	xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx	ccccccccceeeee
25	SEQ SEG PRD MEM	eeeeeccccceeeeeccccccccceeeeeeccccccccc	eeeecccccchhhh
30	SEQ SEG PRD MEM	hhhhhhhhccc	
35	(No	Prosite data available for DKFZphamy2_11d2.3)
	(No	Pfam data available for DKFZphamy2_11d2-3)	

WO 01/98454 DKFZphamy2_11n4

5 group: nucleic acid management

DKFZphamy2_11n4 encodes a novel 1091 amino acid protein with similarity to RAD18 of Schizosaccharomyces pombe and YLR383w of Saccharomyces cerevisia.

The novel protein contains a ATP/GTP-binding site motif A (P-loop). It has similarity to RADLB acts in a DNA repair pathway for removal of UV-induced DNA damage. YLR3B3w of Saccharomyces cerevisiae is a recombination repair protein.

The new protein can find application in modulation of DNA-repair and a as a new tool for manipulation of nucleic acids.

similarity to RADLA (Schizosaccharomyces pombe)

Comment on P53692:
FUNCTION: ACTS IN A DNA REPAIR PATHWAY FOR REMOVAL OF UV-INDUCED DNA DAMAGE THAT IS DISTINCT FROM CLASSICAL NUCLEOTIDE EXCISION REPAIR AND IN REPAIR OF IONIZING RADIATION DAMAGE.

25 Sequenced by EMBL

Locus: /map="2"

30 Insert length: 3679 bp
Poly A stretch at pos. 3646, polyadenylation signal at pos. 3620

L ACCGCGGTGG GCGCCGGGGC TCCCGGGAAT CTACCTTCTC CTGCGGCCGG

35 51 CACGCGGTTC CCAGGGGGCC AGCGGCGGTC AGCCGAGGTC GAGACGCCCG DD CAGGGTGGCC TTAGCGGCCG GTCGTACCAC GGCAGCCCCG CCGATCAGGT 151 TCCTTTGGGA GACTTCGACT TGTTGGCGAA ATGAACCGGA GAAGAATCCC 2D1 AATTGGGAAT TGCGGAAAAC AGGACTCTAG GGTAGAGAAA GGTTGTAGAA 25% CCAATAGGGT TTGAGACCTG ATGGCCAAAA GAAAGGAAGA AAATTTTTCC 40 BOL TCTCCTAAAA ATGCCAAAAG GCCAAGACAA GAAAATTGG AGGATTTTTGA 351 TAAAGATGGT GACGAAGACG AATGTAAAGG TACTACTTTG ACTGCAGCAG 4D% AAGTTGGAAT AATTGAGAGT ATTCACCTAA AAAACTTCAT GTGTCATTCA 45% ATGCTTGGAC CTTTTAAGTT TGGTTCTAAT GTCAACTTTG TTGTTGGCAA 501 CAATGGAAGT GGGAAGAGTG CAGTACTCAC AGCTCTCATA GTCGGTCTTG 551 GTGGAAGAGC AGTTGCTACT AATAGAGGAT CCTCTTTAAA AGGTTTTGTG 45 LOL AAAGATGGAC AGAACTCTGC AGATATCTCA ATAACATTGA GGAACAGAGG LSL AGATGATGCC TTTAAAGCCA GTGTGTATGG TAACTCTATA CTTATACAGC 701 AACACATCAG CATAGATGGA AGTCGATCTT ATAAACTTAA AAGTGCAACA 751 GGCTCCGTGG TTTCCACGAG GAAAGAAGAG CTGATTGCAA TTCTTGATCA BOL TTTTAACATC CAGGTGGATA ATCCAGTTTC TGTTTTAACA CAAGAAATGA 50 B5D GCAAGCAGTT CTTACAGTCT AAAAATGAAG GAGACAAATA CAAATTCTTC **PDL ATGAAAGCAA CGCAACTTGA ACAGATGAAG GAAGATTATT CATACATTAT** 951 GGAAACGAAA GAAAGAACAA AGGAGCAGAT ACATCAAGGA GAAGAGCGGC LOOL TTACTGAACT AAAGCGCCAG TGTGTAGAGA AAGAGGAACG TTTTCAAAGT LD5L ATTGCTGGTT TAAGTACAAT GAAGACTAAT TTAGAGTCCT TGAAACATGA 55 LLDL AATGGCTTGG GCAGTGGTCA ATGAAATTGA AAAACAATTG AATGCCATCA 1151 GAGATAATAT CAAAATTGGA GAAGATCGTG CTGCTAGACT TGACAGGAAA 12D1 ATGGAAGAAC AGCAGGTCAG ACTTAATGAG GCAGAACAAA AGTACAAGGA

1251 TATTCAAGAC AAACTAGAAA AGATTAGTGA AGAGACAAAT GCACGAGCAC TATOODDAA AADAATOTT OTTOTADAA GCAGATGTTG TTGCTAADAA AGGCCCTAT 1351 AATGAAGCTG AGGTTTTATA TAACCGATCC TTAAACGAAT ATAAAGCATT 1401 AAAGAAAGAT GATGAGCAGC TTTGTAAACG AATTGAAGAG CTGAAAAAAA 5 1451 GTACTGACCA ATCTTTGGAA CCTGAACGGT TGGAAAGACA AAAAAAAATA 1501 TCTTGGTTAA AAGAGAGAGT AAAGGCCTTT CAAAATCAAG AAAATTCAGT 1551 CAATCAAGAG ATCGAACAGT TTCAGCAAGC CATAGAAAAG GACAAAGAAG JUD AACATGGCAA AATTAAGAGA GAAGAATTAG ATGTGAAGCA TGCACTGAGC 1651 TACAATCAGA GGCAACTGAA AGAATTGAAA GATAGTAAAA CTGATCGACT 1701 CAAAAGATTT GGCCCTAATG TTCCAGCTCT TCTTGAAGCC ATAGATGATG 10 1751 CTTATAGACA AGGACATTTT ACCTATAAAC CTGTAGGCCC TTTAGGAGCT LABLI TGCATTCATC TTCGGGACCC AGAACTTGCT TTGGCTATTG AATCTTGCTT
LABLI TGCATTCATC TTCGGGACCC AGAACTTGCT TTGGCTATTG AATCTTGCTT
LABLI AAAAGGGCTT CTGCAGGCCT ATTGTTGCCA TAATCATGCT GATGAAAGGG
LABLI TCCTTCAGGC ACTCATGAAA AGGTTTTATT TACCAGGGAC CTCACGGCCA
LABLI CCGATAATAG TTTCTGAGTT TCGGAATGAG ATATATGATG TAAGACACAG
LBLI ATAATGCGGT TGTGGCAAAT AGCCTAATTG ACATGAGAGG CATAGAGACA
LBLI ATAGCCACCA TAAAAATTGTA CACCAGCTT TAGCAGTAA TGCAGTCCAA 15 2301 GTGCTACTAA TCAAAAATAA TTCTGTAGCT CGTGCAGTAA TGCAGTCCCA
2151 AAAGCCACCC AAAAATTGTA GAGAAGCTTT TACTGCTGAT GGTGATCAAG
2201 TTTTTGCAGG ACGTTATTAT TCATCTGAAA ATACAAGACC TAAGTTCCTA
2251 AGCAGAGATG TGGATTCTGA AATAAGTGAC TTGGAGAATG AGGTTGAAAA
2301 TAAGACGGCC CAGATATTAA ATCTTCAGCA ACATTTATCT GCCCTTGAAA
2351 AAGATATTAA ACACAATGAG GAACTTCTTA AAAGGTGCCA ACTACATTAT
2401 AAAGAACTAA AGATGAAAAT AAGAAAAAAT ATTTCTGAAA TTCGGGAACT
2451 TGAGAACATA GAAGAACACC AGTCTGTAGA TATTGCAACT TTGGAAGATG
2501 AAGCTCAGGA AAATAAAAGC AAAATGAAAA TGGTTGAAAA TAGAAGCAGA 20 25 2551 CAACAAAAG AAAATATGGA GCATCTTAAA AGTCTGAAAA TAGAAGCAGA 2601 AAATAAGTAT GATGCAATTA AATTCAAAAT TAATCAACTA TCGGAGCTAG 2651 CAGACCCACT TAAGGATGAA TTAAACCTTG CTGATTCTGA AGTGGATAAC 2701 CAAAAACGAG GGAAACGACA TTATGAAGAA AAACAAAAAG AACACTTGGA 30 2751 TACCTTAAAT AAAAAGAAAC GAGAACTGGA TATGAAAGAG AAAGAACTAG ZBD1 AGGAGAAAT GTCACAAGCA AGACAAATCT GCCCAGAGCG TATAGAAGTA 2851 GAAAAATCTG CATCAATTCT GGACAAGAA ATTAATCGAT TAAGGCAGAA 2901 GATACAGGCA GAACATGCTA GTCATGGAGA TCGAGAGGAA ATAATGAGGC 2951 AGTACCAAGA AGCAAGAGAG ACCTATCTTG ATCTGGATAG TAAAGTGAGG 35 BODL ACTITAAAA AGTITATTAA AATACTGGGA GAAATCATGG AGCACAGTT 3051 CAAGACATAT CAACAATTTA GAAGGTGTTT GACTTTACGA TGCAAATTAT BIDI ACTTTGACAA CTTACTATCT CAGCGGGCCT ATTGTGGAAA AATGAATTTT 3151 GACCACAAGA ATGAAACTCT AAGTATATCA GTTCAGCCTG GAGAAGGAAA 3201 TAAAGCTGCT TTCAATGACA TGAGAGCCTT GTCTGGAGGT GAACGTTCTT 40 3251 TCTCCACAGT GTGTTTTATT CTTTCCCTGT GGTCCATCGC AGAATCTCCT BADD TTCAGATGCC TGGATGAATT TGATGTCTAC ATGGATATGG TTAATAGGAG 3351 AATTGCCATG GACTTGATAC TGAAGATGGC AGATTCCCAG CGTTTTAGAC AAATTATCTT GCTCACACCT CAAAGCATGA GTTCACTTCA CTCAGTAAT 3451 CTGATAAGAA TTCTCCGAAT GTCTGATCCT GAAAGAGGAC AAACTACATT 45 3501 GCCTTTCAGA CCTGTGACTC AAGAAGAAGA TGATGACCAA AGGTGATTTG 3551 TAACTTAACA TGCCTTGTCC TGATGTTGAA GGATTTGTGA AGGGAAAAA AAAAAAAA AAAAAAAA £2dE

PCT/IB01/02050

BLAST Results

55 No BLAST result

50

WO 01/98454

Medline entries

96069417:

Lehmann AR, Walicka M, Griffiths DJ, Murray JM, Watts FZ,

5 McCready S

Carr AM-; The radle gene of Schizosaccharomyces pombe defines a new subgroup of the SMC superfamily involved in DNA repair. Mol Cell Biol 1995 Dec;15(12):7067-80

10 99380167:

Mengiste T. Revenkova E. Bechtold N. Paszkowski J.; An SMC-like protein

is required for efficient homologous

recombination in Arabidopsis. EMBO J 1999 Aug 16:18(16):4505-12

15

Peptide information for frame 1

20

ORF from 271 bp to 3543 bp; peptide length: 1091 Category: similarity to known protein Classification: Nucleic acid management

25 Prosite motifs: RGD (126-128) ATP_GTP_A (76-83)

MAKRKEENFS SPKNAKRPRQ EELEDFDKDG DEDECKGTTL TAAEVGIIES

1 IHLKNFMCHS MLGPFKFGSN VNFVVGNNGS GKSAVLTALI VGLGGRAVAT

101 NRGSSLKGFV KDGQNSADIS ITLRNRGDDA FKASVYGNSI LIQQHISIDG

151 SRSYKLKSAT GSVVSTRKEE LIAILDHFNI QVDNPVSVLT QEMSKQFLQS

201 KNEGDKYKFF MKATQLEQMK EDYSYIMETK ERTKEQIHQG EERLTELKRQ

251 CVEKEERFQS IAGLSTMKTN LESLKHEMAW AVVNEIEKQL NAIRDNIKIG

35 301 EDRAARLDRK MEEQQVRLNE AEQKYKDIQD KLEKISEETN ARAPECMALK

351 ADVVAKKRAY NEAEVLYNRS LNEYKALKKD DEQLCKRIEE LKKSTDQSLE

401 PERLERQKKI SWLKERVKAF QNQENSVNQE IEQFQQAIEK DKEEHGKIKR

451 EELDVKHALS YNQRQLKELK DSKTDRLKRF GPNVPALLEA IDDAYRQGHF

501 TYKPVGPLGA CIHLRDPELA LAIESCLKGL LQAYCCHNHA DERVLQALMK

40 551 RFYLPGTSRP PILVSEFRNE IYDVRHRAAY HPDFPTVLTA LEIDNAVVAN

601 SLIDMRGIET VLLIKNNSVA RAVMQSQKPP KNCREAFTAD GDQVFAGRYY

653 SSENTRPKFL SRDVDSEISD LENEVENKTA QILNLQQHLS ALEKDIKHNE

701 ELLKRCQLHY KELKMKIRKN ISEIRELENI EEHQSVDIAT LEDEAQENKS

751 KMKMVEEHME QQKENMEHLK SLKIEAENKY DAIKFKINQL SELADPLKDE

45 801 LNLADSEVDN QKRGKRHYEE KQKEHLDTLN KKKRELDMKE KELEEKMSQA

852 RQICPERIEV EKSASILDKE INRLRQKIQA EHASHGDREE IMRQYQEARE

901 TYLDLDSKVR TLKKFIKLLG EIMEHRFKTY QQFRRCLTLR CKLYFDNLLS

953 QRAYCGKMNF DHKNETLSIS VQPGEGNKAA FNDMRALSGG ERSFSTVCFI

1001 LSLWSIAESP FRCLDEFDVY MDMVNRRIAM DLILKMADSQ RFRQFILLTP

50 1051 QSMSSLPSSK LIRILRMSDP ERGQTTLPFR PVTQEEDDDQ R

BLASTP hits

55

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_lln4, frame. L

SWISSPROT: RALB_SCHPO DNA REPAIR PROTEIN RADLB., N = L, Score = 1057' b = 5.2e-103

5

PIR:S51470 hypothetical protein YLR383w - yeast (Saccharomyces cerevisiae), N = 1, Score = 823, P = 5e-82

>SWISSPROT: RALB_SCHPO DNA REPAIR PROTEIN RADLB. 10 Length = 1,140

HSPs:

15 Score = 1021 (153.2 bits), Expect = 5.2e-103, P = 5.2e-103Identities = 315/1091 (28%), Positives = 540/1091 (49%)

2 AKRKEENFSSPKNAKRPR@EELEDF--DKDGDEDECKGTTLTAAE----VGIIESIHLKN 55

- 20 AR++N+E ++DG+ D VG+IE IHL N Sbjct: 45 ASRNQDNRPERQSRLQRSSSLIEQVRGNEDGENDVLNQTRETNSNFDNRVGVIECIHLVN 104
- 25 Query: FMCHSMLGPXXXXXXXXXXXXXXXXXXXXXAVLTALIVGLGGRAVATNRGSSLKGFVKDGQN 115 FMCH L A+LT L + LG +A TNR ++K VK G+N

Sbjct: 105 FMCHDSL-

30 KINFGPRINFVIGNGSKSAILTGLTICLGAKASNTNRAPMKSLVKQGKN 163

Query: 114 SADISITLRNRGDDAFKASVYGNSILIQQHISIDGSRSYKLKSATGSVVSTRKEELIAIL 175 A IS+T+ NRG +A++ +YG SI I++ I +GS Y+L+S

35 G+V+ST+++EL Sbict: 164 YARISVTISNRGFEAY@PEIYGKSITIERTIRREGSSEYRLRSFNGTVISTKRDELDNIC 223

Querv: 176

40 DHFNIQVDNPVSVLTQEMSKQFLQSKNEGDKYKFFMKATQLEQMKEDYSYIMETKERTKE 235 DH +Q+DNP+++LTQ+ ++QFL + + +KY+ FMK QL+Q++E+YS I TK Sbjct:

. 224

DHMGLQIDNPMNILTQDTARQFLGNSSPKEKYQLFMKGIQLKQLEENYSLIEQSLINTKN 283

45 Query: 536

QIHQGEERLTELKRQCVEKEERFQSIAGLSTMKTNLESLKHEMAWAVVNEIEKQLNAIRD 295 + ++ L ++ E + ++ LE K EM WA V

E+EK+L

50 Sbict: 284 VLGNKKTGVSYLAKKEEEYKLLWEQSRETENLHNLLEQKKGEMVWAQVVEVEKEL---- 33A

Query: 296 NIKIGEDRAARLDRKMEEQQVRLNEAEQKYKDIQDKLEKISEETNARAP-ECMALKADVV 354

55 + K+ E + L + E DI K+ EE RA E -221X9dIdCsputVIZ3LN3XA3CHQEFQHAEVKLSEAKENLESIVTNQSDIDGKISS-KEEVIGRAKGETDTTKSKFE 395

Query: 355 AKKRAYNEAEVLYNRSLNEYKALKKDDEQLCKRIEELKKSTDQSLEPERLERQKKISWLK 414 Y +N+ K+D + Ι K D Ε ER 5 Sbjct: 396 DIVKTFDG----YRSEMNDVDIQKRDIQN---SINAAKSCLDVYREQLNTERARENNLGG 448 415 ERVKAFQNQENSVNQEIEQF-QQAIEKDKE----EHG----Querv: 10 KIKREELDVKHALS 460 +++ · N+ N++ +EI +Q +E + + E G 2 + + +Sbjct: 449 SQIEKRANESNNLQREIADLSEQIVELESKRNDLHSALLEMGGNLTSLLTKKDSIANKIS 50A 15 Query: 461 YNQRQLKELKDSKTDRLKRFGPNVPALLEAIDDAYRQGHFTYKPVGPLGACIHLRDPELA 520 LK L+D + D++ FG N+P LL+ I R+ · F + P GP+G 20 Sbjct: 509 DQSEHLKVLEDVQRDKVSAFGKNMPQLLKLIT---RETREAHPPKGPMGKYMTVKEQKWH 565 Query: 521. LAIESCLKGLLQAYCCHNHADERVLQALMKRFYLPGTSRPPIIVSEFRNEIYDVRHRAAY 580 25 L IE L ++ + +H D+ +L+ LM++ ++V + Sbjct: 566 LIIERILGNVINGFIVRSHHDQLILKELMRQSNCHAT----VVVGK-----YDPFDYSSG LLL 30 Querv: 581 HPD--FPTVLTALEIDNAVVANSLIDMRGIETVLLIKNNSVARAVRQSQKPPKNCREAFT 638 PD +PTVL ++ D+ V ++LI+ GIE +LLI++ A A M+ 617 EPDSQYPTVLKIIKFDDDEVLHTLINHLGIEKMLLIEDRREAEAYMK--Sbict: 35 RGIANVTQCYA 674 L39 ADG-DQVFAGRYYSSENTR--PKFLSRDVDSEI---SDLENEVENKTAQILNLQQHLSAL 692 D ++ + R S++ + Ι Z EEKL 40 Q + ++ Sbjct: LDPRNRGYGFRIVSTQRSSGISKVTPWNRGTSTSSSSISIEAEKKILDDLKKQYNFASN 734 PA3 E-KDIKHNEELTKKCGTHAKETKWKIKKNIZ-EIKETENIEEHG-ZA-D---Query: 45 IATLEDEA 745 + K KR + Ε I+K I + RE+ ++E + SV DI TLE Sbjct: 735 QLNEAKIEQAKFKRDEQLLVEKIEGIKKRILLKRREVNSLESQELSVLDTEKIQTLERRT 794 50 746 QENKSKMKMVEEHMEQQKENMEH-LKSLKIEAENKYDAIKFKINQLSELADPLKDELN-L 803 E + +++ ++ K N EH ++ +·KI L+ EL+ L 55 Sbjct: 795 SETEKELESYAGQLQDAK-

NEEHRIRDNARPVIEEIRIYREKIATETARLSSLATELSRL 853

PCT/IB01/02050 WO 01/98454

Query: 804

ADSEVDNQKRGKRHYEEKQKEHLDTLNXXXXXXXXXXXXXXXXQARQICPERIEVEKS & & L3 +RH + + + L

ER+ V+ S

5 Sbict: 454 RDEKRNSEVDIERH-RQTVESCTNILREKEAKKVQCAQVVADYTAKANTRC-ERVPVQLS 911

864 ASILDKEINRLRQKIQAEHASHG-Query: DREEIMRQYQEARETYLDLDSKVRTLKKFIKLLGEI 922

10 + LD EI RL+ +I G E+ Y A+E + Sbict: 912

PAELDNEIERL@M@IAEWRNRTGVSVE@AAEDYLNAKEKHD@AKVLVARLT@LL@ALEET 971

15 Query: 923 MEHRFKTYQQFRRCLTLRCKLYFDNLLSQRAYCGKMNFDHKNETLSISVQPGEGNKA-AF 981 + R + + +FR+ +TLR K F+ LSQR + GK+ H+ E L

Sbjct:

20 LRRRNEMUTKFRKLITLRTKELFELYLSQRNFTGKLVIKHQEEFLEPRVYPANRNLATAH 1031

982 N-----DMRALSGGERSFSTVCFILSLWSIAESPFRCLDEFDVYMDMVVRIAMDLIL 1034 ++ LSGGE+SF+T+C +LS+W P RCLDEFDV+MD VNR

25 +++ Sbjct: 1032 NRHEKSKVSVQGLSGGEKSFATICMLLSIWEAMSCPLRCLDEFDVFMDAVNRLVSIKMMV 1091

1035 KMADSQRFRQFILLTPQSMSSLPSSKLIRILRMSDPERGQTTLP 1078 30 +QFI +TPQ M + K + + R+SDP

1092 DSAKDSSDKQFIFITPQDMGQIGLDKDVVVFRLSDPVVSSSALP 1135 Sbjct:

Pedant information for DKFZphamy2_lln4, frame 1

Report for DKFZphamy2_11n4.1

40 ELENGTHE 1091

35

J5P35P-J3 EMMI

6.57 [pI]

SWISSPROT: RALB_SCHPO DNA REPAIR PROTEIN RADLB. Le-[HOMOL] 109

EFUNCATI 03-19 recombination and dna repair ES. cerevisiae, YLR383w3 le-88 EFUNCATE OB-O7 vesicular transport (golgi network, etc.) cerevisiae, YDLO58wl 3e-16 **EFUNCATE** 30.03 organization of cytoplasm ES. cerevisiae,

YDLO58wl 3e-lb 50

> **EFUNCATI** 09.13 biogenesis of chromosome structure EZcerevisiae, YLRO86wl 2e-14

EFUNCATD 1 genome replication, transcription, recombination and EM. jannaschii MJ16431 3e-14 repair

55 30.04 organization of cytoskeleton EFUNCATI ES. cerevisiae, YILl49cl le-l2 03.22 cell cycle control and mitosis [S. cerevisiae. EFUNCATI YDR356wl 8e-12

EFUNCATI 09.10 nuclear biogenesis ES. cerevisiae, YDR356wl 8e-12 **EFUNCATO** 30-10 nuclear organization ES- cerevisiae, YFL008wl 3e-11 EFUNCATI 11.04 dna repair (direct repair, base excision repair ES. cerevisiae, YOR216cll **EFUNCATI** 99 unclassified proteins 5e-09 myosin-l isoforml &e-O& 10 **LFUNCATI** 03.04 budding, cell polarity and filament formation ES. cerevisiae, YHRO23W MYOJ - myosin-l isoforml &e-O& EFUNCATI DB-22 cytoskeleton-dependent transport ES- cerevisiae-YHRD23w MYOL - myosin-l isoforml &e-O8 EFUNCATI Ob.O7 protein modification (glycolsylation, acylation, 15 myristylation, palmitylation, farnesylation and processing)
ES. cerevisiae, YKL201c1 2e-07 EFUNCATI D3.13 meiosis ES. cerevisiae, YDR285wl 4e-D7 30.13 organization of chromosome structure **EFUNCATI** cerevisiae, YDR285wJ 4e-07 20 EFUNCATI 98 classification not yet clear-cut ES. cerevisiae. YJR134c1 7e-07 EFUNCATI Ob-10 assembly of protein complexes ES. cerevisiae. YPR141c1 7e-07 25 **EFUNCATI** 30.05 organization of centrosome ES. cerevisiae. YPR141c3 7e-07 ll.Dl stress response ES. cerevisiae, YPR141cl 7e-07 **EFUNCATE** EFUNCATI 03.07 pheromone response, mating-type determination, 30 [FUNCAT] r general function prediction TH. influenzae, HI07561 le-06 EFUNCATI 10-05-99 other pheromone response activities EZcerevisiae, YHR158cJ 2e-06 EFUNCATE 05.04 translation (initiation, elongation and ES. cerevisiae, YALO35wl 3e-04 35 termination) **EFUNCATE** 30-02 organization of plasma membrane ES. cerevisiae. YER008c3 4e-04 **EFUNCATI** 08.16 extracellular transport ES. cerevisiae. YERDD&c3 4e-04 40 **EFUNCATI** 09.04 biogenesis of cytoskeleton ES cerevisiae. YKL179c1 7e-04 [FUNCAT] 03.22.01 cell cycle check point proteins EZcerevisiae, YGLO86w3 7e-04 [FUNCAT] 08.01 nuclear transport [S. cerevisiae, YDL207w] [].001 ES. cerevisiae, YDL207wl 0.001 45 EFUNCATI 04-07 rna transport BL00326C Tropomyosins proteins **EBFOCK21** PRO1004B [BLOCK2] EBFOCK21 BLOOL21A Colipase proteins **EBFOCK21** PF00580A d2tmab_ 1.105.4.1.1 Tropomyosin Erabbit 50 EZCOPI (Oryctolagus cuniculus) 3e-06 3.6.1.32 Myosin ATPase 9e-20 TEC]

EPIRKUJ

phosphotransferase 9e-16

PCT/IB01/02050 endocytosis 2e-13

WO 01/98454

EPIRKWI heart 9e-2D **EPIRKWI EPIRKWI** polymorphism le-10 **EPIRKWI** serine/threonine-specific protein kinase 9e-16 5 **CPIRKWI** transmembrane protein &e-15 **TPIRKW** zinc finger 2e-13 metal binding 2e-13 **EPIRKW**3 **CPIRKWI** DNA binding 2e-06 **EPIRKUB** muscle contraction 9e-20 10 **EPIRKUI** acetylated amino end 3e-13 **EPIRKU**J actin binding 9e-20 **EPIRKU** mitosis &e-10 **EPIRKUI** microtubule binding 3e-09 **EPIRKUI** chromosomal protein 3e-11 15 EPIRKWI · ATP 9e-20 **TPIRKUB** receptor 2e-06 thick filament 9e-20 **EPIRKUI EPIRKU** phosphoprotein 2e-14 **EPIRKUD** glycoprotein le-lD 20 **EPIRKWI** skeletal muscle le-l8 calcium binding Ze-10 **EPIRKU**I **EPIRKUI** alternative splicing 3e-12 **EPIRKW3** DNA condensation 3e-11 **EPIRKWI** P-loop 9e-20. 25 **EPIRKUB** coiled coil 9e-20 **EPIRKWI** heptad repeat le-10 methylated amino acid 9e-20 **CPIRKWI CPIRKUI** basement membrane le-10 **EPIRKW3** immunoglobulin receptor 4e-09 peripheral membrane protein 2e-13 30 **CPIRKWI CPIRKU**I cardiac muscle 9e-20 **CPIRKUI** extracellular matrix le-10 **EPIRKU**I hydrolase 9e-20 **CPIRKWI** microtubule 2e-10 35 **EPIRKUI** muscle 2e-14 **EPIRKUJ** membrane protein le-10 **EPIRKUI** EF hand 2e-10 **EPIRKUI** cell division &e-10 **EPIRKUI** cytoskeleton le-13 40 **EPIRKWI** hair 2e-10 **EPIRKWI** calmodulin binding 2e-13 **EPIRKWI** Golgi apparatus Le-O8 **IPIRKUJ** smooth muscle 2e-07 **ESUPFAMI** conserved hypothetical P115 protein 4e-26 45 ESUPFAMI myosin heavy chain 9e-20 ESUPFAMI unassigned Ser/Thr or Tyr-specific protein kinases 9e-JP **EZUPFAMJ** centromere protein E 3e-09 **ESUPFAME** calmodulin repeat homology 2e-10 50 **EZUPFAMD** alpha-actinin actin-binding domain homology 7e-07 EZUPFAMI myosin motor domain homology 9e-20 **ESUPFAMD** tropomyosin 5e-08 ESUPFAMD plectin 7e-07 ESUPFAME pleckstrin repeat homology 3e-09 55 **ESUPFAMI** trichohyalin 2e-10 **ESUPFAMI** hypothetical protein MJ1322 2e-06 ESUPFAMI ribosomal protein SLO homology 7e-07

	wo	01/98454						F	CT/IB01/020)50	
	EZUP	FAMB	giantin								
		PFAMB	protein	kinase	homolo	ogy 9e-	-1P				
		PFAMI		motor d							
5		PFAMB PFAMB		arly end ein 4e-0		antige	su T Se.	- т ⁻			
3		PFAMI		letal ke		Ao-NL					
			ATP_GTP		deli	06 00					
			RGD 1				•				
	EKWI	_	All_Alp	ha							
10	EKUI			PLEXITY							
	EKWI	J	COILED_	COIL	7 P	15 %			•		
			•								
	SEQ	MAKR	KEENFSSP	KNAKRPRQ	EELEDF	DKDGDE	DECKGT	TLTAAE	GIIESIH	LKNFMC	21
15	SEG										• •
	PRD COIL		hhhhcccc	ccccchh	hhhhc	ccccc	ccccc	cccccc	ceeeeeh	hhhhhc	CC
	COIL										
					,						• •
20	SEQ		FKFGSNVN								
	SEG		×××××××								
	PRD COIL		cccceee	eeecccc	ccnnni	nnnnnn	iccccc	cccccc	ceeeecc	cccee	2 e
	CVIL										
25											
	SEQ		NRGDDAFK							AILDHF	NI
	SEG PRD		ccccccc							 hhhhhhhi	 hh
	COIL				CCIIIIII	mecece	.eeeecc	ccciiiiii			. 11 1
30		• • • •						• • • • • •		• • • • • •	• •
	SEQ	AUDN	PVSVLTQE	MCKVEI VC	NECNI	VONCEMA	/ A T / L T / L	MVERUEL	ノエMピサレビロ	TVEATU	^
	SEG	K A D IA	PAZAFIKE		NIVE GDI		AIRLER		THEIVER	IVERTH	טא
	PRD	cccc	hhhhhhhh	hhhhhhhh	hhcchl	hhhhhh	hhhhhh	hhhhhhh	hhhhhhh	hhhhhhl	hh
35	COIL	Ζ.									
			• • • • • • • •	• • • • • • • •	• • • • • •	• • • • • •		 .	• • • • • • •	• • • • • •	• •
	SEQ	EERL'	TELKRQCV	EKEERF@S:	IAGLS	IMKTNLE	SLKHEM	AWAVVNE	IEKQLNA	IRDNIK	IG
	SEG										• •
40	PRD	_	hhhhhhhhh	hhhhhhhhl	hhhhhh	hhhhhhh	ւիհիիիի	hhhhhhh	hhhhhhhh	hhhhhhl	nh
	COIL	. Z						,		CCCCCC	<i></i>
		••••			• • • • •			• • • • • • • •			د ر
	SEQ	EDRA	ARLDRKME	EQQVRLNE	AEQKYK	(DIQDKL	EKISEE	TNARAPE	CMALKAD	VVAKKR	ΑY
45	SEG		• • • • • • • • • • • • • • • • • • •					• • • • • •		· • • • • • •	• •
	PRD COIL		hhhhhhhhh	hnnhhnnh	nnnnn	որորոր	ւրրրիրի	hhhhhhh	ւրրրրիր	hhhhhhl	nh
	COIL		ccccccc	ccccccc	ccccc	ccccc	ccccc				
50	SEQ	NEAE	VLYNRSLNI					•		LKERVK	AF
	SEG	 In In In In In I									• •
	PRD COIL		hhhhhhhhl	เเกก กกกกก)	unnnh	innn nh	innnnhh.	nnnnnn	เกิดกาทที่ที่	nnnnhh	าท
	CATE					. .	· • • • • • •		· • • • • • • •		
55					• *	•					
	SEQ	QNQE	NZVNGEIE							KTDRLKI	RF
	SEG PRD	hbbb!	 hhhhhhhhhl								• • h.h.
	rk n	mmini			411111111		manann.	nananar	เกษา	ומממחוז	111



		••••••
5	SEQ. SEG PRD	GPNVPALLEAIDDAYRQGHFTYKPVGPLGACIHLRDPELALAIESCLKGLLQAYCCHNHA
•	COIL	2
10	SEQ SEG	DERVLQALMKRFYLPGTSRPPIIVSEFRNEIYDVRHRAAYHPDFPTVLTALEIDNAVVAN
	PRD COIL:	hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
15	SEQ	
	SEG PRD	SLIDMRGIETVLLIKNNSVARAVMQSQKPPKNCREAFTADGDQVFAGRYYSSENTRPKFL
20	COIL	
	SEQ	
	SEG	SRDVDSEISDLENEVENKTAQILNLQQHLSALEKDIKHNEELLKRCQLHYKELKMKIRKN
25	PRD COILS	
		CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
	SEQ SEG	ISEIRELENIEEHQSVDIATLEDEAQENKSKMKMVEEHMEQQKENMEHLKSLKIEAENKY
30	PRD COILS	հիրորդ անական անական հետում և հետում և Տ
35	SEG SEG	DAIKFKINGLSELADPLKDELNLADSEVDNQKRGKRHYEEKQKEHLDTLNKKKRELDMKE
	PRD COILS	
	•	CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
40	SEG	KELEEKMS@AR@ICPERIEVEKSASILDKEINRLR@KI@AEHASHGDREEIMR@Y@EARE
	PRD COILS	hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
45		cccccccccc
	SEG SEG	TYLDLDSKVRTLKKFIKLLGEIMEHRFKTYQQFRRCLTLRCKLYFDNLLSQRAYCGKMNF
50	PRD COILS	hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
50		
	SEQ SEQ	DHKNETLSISVQPGEGNKAAFNDMRALSGGERSFSTVCFILSLWSIAESPFRCLDEFDVY
55	PRD COILS	eccccceeeecccchhhhhhccccccchhhhhhhhhhhh
	SEQ	MDMVNRRIAMDLILKMADS@RFR@FILLTP@SMSSLPSSKLIRILRMSDPERG@TTLPFR
		TO THE PERSON OF

	WO 01/98454	•		PCT/IB01/02050
	SEGPRD hhhhhh COILS	րիրերեր և հեռուս և հ		ccceeeeeccccccccccccccccccccccccccccc
5	SEQ PVTQEEI		•••••	• • • • • • • • • • • • • • • • • • • •
10	PRD chhhhh	hccc		
		Prosit	e for DKFZphamyi	2_lln4.l
15	6200075 620007P	126->129 R <i>G</i> 76->84 AT	D P_GTP_A	PD0C00016 PD0C00017
20	(No Pfam dat	a available for	DKFZphamy2_lln	4-3)

DKFZphamy2_121f19

5 group: cell structure and motility

DKFZphamy2_121f19 encodes a novel 251 amino acid protein with high similarity to a Rat ankyrin binding glycoprotein-1 related mRNA.

Ankyrin binding glycoproteins play a role in neural cell adhesion and in prosate tumor cell transformation. DKFZphamy2_121f19.p3 is expressed in brain, uterus and prostate above average.

- 15 The new protein can find application modulation of cyto skeletonmembrane interactions.
- similarity to ankyrin binding glycoprotein-1 related mRNA (Rattus 20 norvegicus)

Sequenced by DKFZ

Locus: /map="l"

25

Insert length: 1498 bp

Poly A stretch at pos. 1479, polyadenylation signal at pos. 1460

1 CGGCACCTTC GCCGGCGCCC TCGCCCACCC CAGCCCCGCC CCAGAAGGAG 30 51 CAGCCCCCG CGGAGACCCC TACAGACGCT GCTGTCTTGA CCTCACCCCC 101 AGCCCCTGCT CCCCCGGTGA CCCCTAGCAA ACCAATGGCC GGCACCACAG 151 ACCGAGAAGA AGCCACTCGG CTCTTGGCTG AGAAGCGGCG CCAGGCCCGG 201 GAGCAGCGGG AGCGCGAGGA GCAGGAGCGG AGGCTGCAGG CAGAAAGGGA 35 251 CAAGCGAATG CGAGAGGAGC AGCTGGCACG GGAGGCCGAG GCCCGGGCGG BAAGAGAGC GGAGGCCCGG AGGCGGGAGG AGCAGGAGGC ACGAGAAGG 351 GCGCAGGCCG AGCAGGAGGA GCAGGAGCGG CTGCAGAAGC AGAAAGAGGA 4D1 GGCCGAAGCT CGGTCGCGGG AAGAGGCGGA GCGGCAGCGT CTGGAGCGGG 451 AAAAGCACTT CCAGCAGCAG GAGCAAGAGC GGCAAGAGCG CAGAAAGCGT 40 501 CTGGAGGAGA TCATGAAGAG GACTCGGAAG TCAGAAGTTT CTGAAACCAA 551 GAAGCAGGAC AGCAAGGAGG CCAACGCCAA CGGTTCCAGC CCAGAGCCTG LOD TGAAAGCTGT GGAGGCTCGG TCCCCAGGGC TGCAGAAGGA GGCTGTGCAG LSI AAAGAGGAGC CCATCCCACA GGAGCCTCAG TGGAGTCTCC CAAGCAAGGA 701 GTTGCCAGCG TCCCTGGTGA ATGGCCTGCA GCCTCTCCCA GCACACCAGG 45 751 AGAATGGCTT CTCCACCAAC GGACCCTCTG GGGACAAGAG TCTGAGCCGA BD1 ACACCAGAGA CACTCCTGCC CTTTGCAGAG GCAGAAGCCT TCCTCAAGAA 851 AGCTGTGGTG CAGTCCCCGC AGGTCACAGA AGTCCTTTAA GAGGGTTTGC 901 CTTGGATCCG GGCACAGTTG TGAGGGCTCC TCTGCATCAC CTACCAGGAT 951 GTCTGGAGGA GAAAAAGACA GAACAAAGAT GGAAGTGGCC TGGGCCCCTG 1001 GGGGTGGGTC CTCTCTGTTG TTTTTAATCT GCACCTTATA GACTGATGTC 50 1051 TCTTTGGCCG GAGCCAGATC TGCCCCTCAG TGCATTCGTG TGCTCGCACG 1101 CGCAGACATC CCTTCTCCCC CATACACACA TATACACTCA CAGCCTCTCT 1151 GGCCTCTTCC CTTGGGGAGG GGCCACCTGT AGTATTTGCC TTGATTTGGT 1201 GGGGTACAGT GGATGTGAAT ACTGTAAATA GCTTGTGCTC AGACTCCTCT 55 1251 GCGTGGAGAG GGTGGGTGCA GGAGGCAGAC CCTCCCCCA AAGCCCCCTG 1301 GGGAGATCTT CCTCTCTCTA TTTAACTGTA ACTGAGGGGG ATCCCAGGTC 1351 TGGGGATGGG GGACACCTTG GGCCACAGGA TACTGGTTGC TTCAGGGGTA 1401 CCCATGCCCC CTGCCCTCGC CTGGAATCAG TGTTACTGCA TCTGATTAAA

1451 TGTCTCCAGA AATAAAGAAT AATTCTGCCA AAAAAAAAA AAAAAAAA

BLAST Results

No BLAST result

10 Medline entries

No Medline entry

15

5

Peptide information for frame 3

20 ORF from 135 bp to 887 bp; peptide length: 251 Category: putative protein Classification: Cell signaling/communication

1 MAGTTDREEA TRLLAEKRRQ AREQREREEQ ERRLQAERDK RMREEQLARE 51 AEARAEREAE ARRREEGEAR EKAQAEGEEQ ERLQKQKEEA EARSREEAER 25 101 QRLEREKHFQ QQEQERQERR KRLEEIMKRT RKSEVSETKK QDSKEANANG 151 SSPEPVKAVE ARSPGLQKEA VQKEEPIPQE PQWSLPSKEL PASLVNGLQP 201 LPAHQENGFS TNGPSGDKSL SRTPETLLPF AEAEAFLKKA VVQSPQVTEV 251 L

30

BLASTP hits

35 No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_121f19, frame 3

No Alert BLASTP hits found

40

Pedant information for DKFZphamy2_121f19, frame 3

Report for DKFZphamy2_121f19-3

45

CLENGTHD 295 EMMI 33517.96 [[q]

TREMBLNEW: AB033013_1 gene: "KIAA1187"; product: 50 EHOMOLI "KIAAll&7 protein"; Homo sapiens mRNA for KIAAll&7 protein, partial cds. Le-64 EBLOCKZI PFO1140D

EBFOCK23

BLDD412D Neuromodulin (GAP-43) proteins

55 **EBFOCK21** BLOD85PC

EBFOCK23 BLDD422C Granins proteins

EBFOCKZ PR00167C

A5PP0079 **EBFOCK21**



	WO 01/984	54	•				PCT/IB01/020	50
	EBFOCKZ:	PR00049D PR00910A	Clathrin	light	chain	proteins		
5	EKM] EKM] EKM]	All_Alpha Low_compl Colleb_co	EXITY	51.19 10.51			·	
10	SEG xx	AGGATGZGAGZ	:ccccccc	cccee	(XXXXX)	CCCCCCC	×××	hhhhhhhhh
15	SEG xx: PRD hhl	R@ARE@REREE@ xxxxxxxxxx nhhhhhhhhhhhh	xxxxxxxx hhhhhhhhhh	hhhhhhl	xxxxxx nhhhhh	(XXXXXXXX hhhhhhhhh	xxxxxxxx hhhhhhhhhh	xxxxxxx hhhhhhhh
20 25	SEQ QEI SEG XXI PRD hhl	EQERLQKQKEE/ xxxxxxxxxx hhhhhhhhhhhh CCCCCCCCCCC	NEARSREEAE (XXXXXXX) Iddddddddddddd	ERQRLEI ×××××× 1hhhhhl	REKHF@G xxxxxxx nhhhhh	QEQERQER XXXXXXXX hhhhhhhhh	RKRLEEIMK xxxxxx hhhhhhhhhh	RTRKSEVS
30	SEQ ET	KKQDSKEANANO hhhhhhhhhccco	SSPEPVKA	VEARSP(SL&KEAY Cchhhhl	V&KEEPIP hhhhcccc	EP@M2Tb2K	CCCCEEEE
35	SEG	QPLPAHQENGFS	ccccccc		chhhhhl	որների Մարդերի		ccc
40		site data av m data avail						

5 group: cell cycle

DKFZphamy2_121m2 encodes a novel 480 amino acid protein with similarity to human PA26-T2 protein.

- 10 PA26-T2 is a p53 responsive gene. The protein is predominantly expressed in brain, breast and kidney and may represent a potential novel regulator of cellular growth. Isoforms are differentially induced by genotoxic stress (UV, gamma-irradiation and cytotoxic drugs) in a p53-dependent manner.
 - The new protein can find application in modulating cell division and apoptosis pathways.
- 20 similarity to PA26 nuclear protein isoforms (Homo sapiens) probably differential polyadenylation

Sequenced by DKFZ

25 Locus: unknown

15

30

Insert length: 3327 bp
Poly A stretch at pos. 3306, polyadenylation signal at pos. 3279

L TCCAGCACCA AAGCGGCCGT TCTCGGATTC CGGAGCGTTC TGGAGCCCCG 51 AGAGACGCCC CGGGGTTCTA GAAGCTCCCC GGCGGCGCCC AGTCCCGGCT LOL TCATTCGGGC GTCCCTCCGA AACCCACTCG GGTGCACGGG TCGTCGGCGA 151 GCCGCGACCG GGTCCTGGCG CGCACCATGA TCGTGGCGGA CTCCGAGTGC 35 201 CGCGCAGAGC TCAAGGACTA CCTGCGGTTC GCCCCGGGCG GCGTCGGCGA 251 CTCGGGCCCC GGAGAGGAGC AGAGGGAGAG CCGGGCTCGG CGAGGCCCTC TODDDDDADD DOTTOTODA DARDTDOTO OTCOTTOTO GASGOCT 351 GAGAGCCTCG AGCAGCACCT GGGGCTGGAG GCACTGATGT CCTCTGGGCG 401 AGTAGACAAC CTGGCAGTGG TGATGGGCCT GCACCCTGAC TACTTTACCA 40 451 GCTTCTGGCG CCTGCACTAC CTGCTGCTGC ACACGGATGG TCCCTTGGCC 501 AGCTCCTGGC GCCACTACAT TGCCATCATG GCTGCCGCCC GCCATCAGTG 55% TTCTTACCTG GTAGGCTCCC ACATGGCCGA GTTTCTGCAG ACTGGTGGTG LOT ACCCTGAGTG GCTGCTGGGC CTCCACCGGG CCCCCGAGAA GCTGCGCAAA L51 CTCAGCGAGA TCAACAAGTT GCTGGCGCAT CGGCCATGGC TCATCACCAA 45 7DL GGAACACATC CAGGCCTTGC TGAAGACCGG CGAGCACACT TGGTCCCTGG 75% CCGAGCTCAT TCAGGCTCTG GTCCTGCTCA CCCACTGCCA CTCGCTCTCC BOD TCCTTCGTGT TTGGCTGTGG CATCCTCCCT GAGGGGGGATG CAGATGGCAG B51 CCCTGCCCC CAGGCACCTA CACCCCCTAG TGAACAGAGC AGCCCCCCAA POL GCAGGGACCC GTTGAACAAC TCTGGGGGCT TTGAGTCTGC CCGCGACGTG 50 951 GAGGCGCTGA TGGAGCGCAT GCAGCAGCTG CAGGAGAGCC TGCTGCGGGA 1001 TGAGGGGACG TCCCAGGAGG AGATGGAGAG CCGCTTTGAG CTGGAGAAGT LOST CAGAGAGCCT GCTGGTGACC CCCTCAGCTG ACATCCTGGA GCCCTCTCCA LLDL CACCCAGACA TGCTGTGCTT TGTGGAAGAC CCTACTTTCG GATATGAGGA
LLSL CTTCACTCGG AGAGGGGCTC AGGCACCCCC TACCTTCCGG GCCCAGGATT 1201 ATACCTGGGA AGACCATGGC TACTCGCTGA TCCAGCGGCT TTACCCTGAG
1251 GGTGGGCAGC TGCTGGATGA GAAGTTCCAG GCAGCCTATA GCCTCACCTA
1301 CAATACCATC GCCATGCACA GTGGTGTGGA CACCTCCGTG CTCCGCAGGG

WO 01/98454



	WO 01/	98454			T P	CT/IB01/02050
	1351	CCATCTGGAA	CTATATCCAC	TGCGTCTTTG	GCATCAGATA	TGATGACTAT
	1401	GATTATGGGG	AGGTGAACCA		CGGAACCTCA	AGGTCTATAT
	1451	CAAGACAGTG	GCCTGCTACC	CAGAGAAGAC	CACCCGAAGA	ATGTACAACC
	1501	TCTTCTGGAG	GCACTTCCGC		AGGTCCACGT	GAACTTGCTG
5	1551	CTCCTGGAGG	CGCGCATGCA	AGCCGCTCTG	CTGTACGCCC	TCCGTGCCAT
	1601	CACCCGCTAC	ATGACCTGAC	TCCTGAGCAG	GACCTGGGCC	CGGTTCAGCT
	1651	CCCCACAAGG	ACTTCTCTGT	CTGGAGACAG	CCCCAGACCC	TTTTGTGTCC
	1701	CATGCCCACC	CTCCCCACGC	TGCAGTGGGC	TTGTGTGTGA	TGTGCAGTCC
	1751	CGAAGCCACA	CCCTCCCTTT	TCCTCACTGG	AATGGACAGT	TCATTGCACT
10	1801	GACTCTGGGA	TCTCAGCCCT	GCTCCTGGGA	GCTGGAAGAG	CACTTGGAGA
	1851	TCCTAAGGGA	CCACACCCTT	CCTCCTTCCC	CTGCCCACAG	AGGCAGAGGG
	1901	CACAGGAAAG	AAGCCGGGCC	AAGCTCGGAA	TTAATGTGCC	ACAAGTGTTG
	1951	TGGCCTTCCT	GAACTGGGAA	GTCCCTGGCT	GGCCCCCGGG	GGAGAGGGGC
	5007	AAATGCCTCC	GGGACTGACA	CTCCAGGCAG	CTTTGCCTTC	TCTCCCCTGT
15	2051	CATTTCCAGA	TTTCATTACC	TCCTACTTGC	CATTCACCCA	TCAATGTGAA
	5707	AGTCAGGGTC	ACAGCTGGTC	TGTGTGTCCA	GTTCCCTAAA	AGCCTGTTCT
	5727	GTTGGGCAGC	CTGAGGCTGT	TGCCCGAATC	CTAGTTCAGT	TTTTTGACTT
	5507	CCTTTGCCCT	TTTTCCCTTT	TCTCCATGCT	TAATGGTGTG	AGGCGTCAGG
	225I	AGAGAGGCCA	AGTACATAAA	AAAAAAAAA	AGCAGATTAT	CTCTAGAGAG
20	5307	TTTGAGCCTT	TGCTGGTCAC	ATTGCCTTCT	GAAGAGGAGG	GAGTATTAGA
	2351	TTATAAATCC	TCTTTATTTT	GGTCCTTTAT	GCTTGAGGTT	CCAACCTGGA
	2401	GCCACAGTGT	GTGAGAGGAG	GAGGAGAGGG	AGAATTCTGT	TCTCCCAGAG
	2451	CTGCACCTGC	CTCGCAGAGG	CCAGCACCCC	ACTCTCCTGC	CTCCAGTGGC
	2501	CCTGCCGCAG	ATGTCTCCCA	AAAAGTTGAG	CCTTTCTAGA	TGGCTTAGGT
25	2551	GGCACCATGG	CTCAGCAGGA	GGGGCGGGAG	GCACCAGGGT	TCTTGTTTGG
	SPOT	ACCCTGCCC	TGGGCCATGG	CCAGGTGACC	ATGGCTACAT	TGCCAAACCT
	5627	CTGACTGCCA	CAGCTGCAGA	CTGAGAGGGT	GGGTCTGAGT	CCCCACAATG
	2701	TCTGAAGCTG	CCCCTGGGAT	TCTCAGGCCA	ACCTGCCAAC	AGCAAGCGGA
	2751	TTTTCTTGCA		CCCCATTTCT	GCAGCCAGTG	TCTCCTGGGT
30 .	5907	GCCTTCTGAG	GACTCCCACC	CCCATCCCAG	TATCTCATCT	GTCCCCTCTC
	2851	CTGGGGCTTA	AGTGGGTTGC	TTCCAGGCAG	AAGCAGCCAA	GGACCGATTC
	2901	CAGGCACTTT	CTGTAGCAAA	TGACTGTGAA	TTACGACTTC	TCTTGCCCTT
	2951	CTTCTAGCAG	TCTGTGCCTC	CTCTCTGACC	AGTTTGGAGG	GCACTGAAGA
25	3007	AAGGCAAGGG	CCGTGCTGCT	GCTGGGCGGG	GCAGGAGAGG	AGCCTGGCCA
35	3051	GTGTGCCACA	TTAAATACCC		GGAGAAGCAA	CCGGCACCCC
	3707		GAAAGCCCTC		GGTGTGCAGG	AGAGAAGAGG
	3151	CCCCGGCATG	GGGATCTGGG	TTCTAGAGGG	CATGTGATGA	CTGTAAATGT
	3527 3507	TCACTGGGTG GGCTTTGCTA	GGTAGGGAGT CCAGTTCCAT	GGTATCCAGT	GTTCAAGTGC ATAAACGTTC	AGAAATCTTT
40	3307	TGTTTCATAA	AAAAAAAAA		ATAAACGTIC	GCTGAGGTTT
40	ココロカ	IGITICATAA	AAAAAAAAAA	AAAAAA		

BLAST Results

45 No BLAST result

Medline entries

•

50

95024170:
Buckbinder L., Talbott R., Seizinger B.R., Kley N.; Gene regulation by

55 temperature-sensitive p53 mutants: identification of p53 response genes. Proc. Natl. Acad. Sci. U.S.A. 91(22):10640-10644(1994).

9124117:

Velasco-Miguel S. Buckbinder L. Jean P. Gelbert L. Talbott R. Laidlaw

J₁ Seizinger B₁ Kley N₂ PA2b₁ a novel target of the p53 tumor suppressor and member of the GADD

5 family of DNA damage and growth arrest inducible genes. Oncogene 1999
Jan 7:18(1):127-37

10

Peptide information for frame 3

15 ORF from 177 bp to 1616 bp; peptide length: 480 Category: strong similarity to known protein Classification: Cell division

451 EKVHVNLLLL EARMQAALLY ALRAITRYMT

1 MIVADSECRA ELKDYLRFAP GGVGDSGPGE EQRESRARRG PRGPSAFIPV

51 EEVLREGAES LEQHLGLEAL MSSGRVDNLA VVMGLHPDYF TSFURLHYLL

101 LHTDGPLASS WRHYIAIMAA ARHQCSYLVG SHMAEFLQTG GDPEWLLGLH

151 RAPEKLRKLS EINKLLAHRP WLITKEHIQA LLKTGEHTWS LAELIQALVL

201 LTHCHSLSSF VFGCGILPEG DADGSPAPQA PTPPSEQSSP PSRDPLNNSG

251 GFESARDVEA LMERMQQLQE SLLRDEGTSQ EEMESRFELE KSESLLVTPS

25 301 ADILEPSPHP DMLCFVEDPT FGYEDFTRRG AQAPPTFRAQ DYTWEDHGYS

351 LIQRLYPEGG QLLDEKFQAA YSLTYNTIAM HSGVDTSVLR RAIWNYIHCV

401 FGIRYDDYDY GEVNQLLERN LKVYIKTVAC YPEKTTRRMY NLFWRHFRHS

30

BLASTP hits

No BLASTP hits available

35

Alert BLASTP hits for DKFZphamy2_121m2, frame 3

TREMBL:AF033120_1 gene: "PA26"; product: "p53 regulated PA26-T2 nuclear

- 40 protein": Homo sapiens p53 regulated PA26-T2 nuclear protein (PA26)
 mRNA: complete cds: N = 1: Score = 1377: P = 9:7e-141
- TREMBL:AFO33122_1 gene: "PA26"; product: "non-p53 regulated PA26-45 T1 nuclear protein"; Homo sapiens non-p53 regulated PA26-T1 nuclear protein (PA26) mRNA; complete cds.; N = 1; Score = 1363; P = 3e-139
- 50 TREMBL:AF033121_1 gene: "PA26"; product: "p53 regulated PA26-T3 nuclear protein"; Homo sapiens p53 regulated PA26-T3 nuclear protein (PA26)

 mRNA, complete cds., N = 1, Score = 1307, P = 2.5e-133

55

>TREMBL:AF033120_1 gene: "PA26"; product: "p53 regulated PA26-T2 nuclear

protein": Homo sapiens p53 regulated PA26-T2 nuclear protein (PA26) mRNA: complete cds:

Length = 492

5 HSPs:

Score = 1377 (206.6 bits), Expect = 9.7e-141, P = 9.7e-141 Identities = 277/471 (58%), Positives = 334/471 (70%)

10
Query: 22 GVGDSGPGEEQRESRARRGPR----GPSAFIPVEEVLREGAESLEQHLGLEALMSSGRV 76

G GG+QE RPR GPS FIP +E+L+ G+E + H L

++ + GR+

15 Sbjct: 22 GCKQCGGRDQDEELGIRIPRPLGQGPSRFIPEKEILQVGSEDAQMHALFADSFAALGRL 81

Query: 77
DNLAVVMGLHPDYFTSFWRLHYLLLHTDGPLASSWRHYIAIMAAARHQCSYLVGSHMAEF 136
DN+ +VH H Y SF + + LL DGPL +RHYI

20 DN+ +VM HP Y SF + + LL DGPL +RHY:
IMAAARHQCSYLV H+ +F
Sbjct: 82

DNITLVMVFHPQYLESFLKTQHYLLQMDGPLPLHYRHYIGIMAAARHQCSYLVNLHVNDF 141

25 Query: 137
LQTGGDPEWLLGLHRAPEKLRKLSEINKLLAHRPWLITKEHIQALLKTGEHTWSLAELIQ 196
L GGDP+WL GL AP+KL+ L E+NK+LAHRPWLITKEHI+ LLK
EH+WSLAEL+
Sbict: 142

30 LHVGGDPKWLNGLENAPQKLQNLGELNKVLAHRPWLITKEHIEGLLKAEEHSWSLAELVH 201

Query: 197 ALVLLTHCHSLSSFVFGCGILPEGDADGXXXXXXXXXXX-----XXXXXXXXXXDPLNNS 249
A+VLLTH HSL+SF FGCGI PE DG

35 P+N++

Sbjct: 202
AVVLLTHYHSLASFTFGCGISPEIHCDGGHTFRPPSVSNYCICDITNGNHSVDEMPVNSA 261

Query: 250 GGF---ESARDVEALMERMQQLQESLLRDEG-

40 TSGEEMESRFELEKSESLLVTPSADILE 305

KF

+S +VEALME+M+QLQE RDE SQEEM SRFE+EK ES+ V S+D E

Sbjct: 262 ENVSVSDSFFEVEALMEKMRQLQEC-RDEEEASQEEMASRFEIEKRESMFVF-SSDDEE 318

45 Querv: 306

55

PSPHPDMLCFVEDPTFGYEDFTRRGAQAPPTFRAQDYTWEDHGYSLIQRLYPEGGQLLDE 365
+P + ED ++GY+DF+R G P TFR QDY WEDHGYSL+

RLYP+ GQL+DE

50 Sbjct: 319 VTPARAVSRHFEDTSYGYKDFSRHGHHVP-TFRVQDYCWEDHGYSLNRVJSVGVEDE 377

AY+LTYNT+AMH

Query: 366 KFQAAYSLTYNTIAMHSGVDTSVLRRAIWNYIHCVFGIRYDDYDYGEVNQLLERNLKVYI 425

VDTS+LRRAIWNYIHC+FGIRYDDYDYGE+NQLL+R+ KVYI Sbjct: 378 KFHIAYNLTYNTMAMHKDVDTSMLRRAIWNYIHCMFGIRYDDYDYGEINQLLDRSFKVYI 437

Query: 426

KTVACYPEKTTRRMYNLFWRHFRHSEKVHVNLLLLEARMQAALLYALRAITRYMT 480 KTV C PEK T+RMY+ FWR F+HSEKVHVNLLL+EARMQA

5 LLYALRAITRYMT Sbjct: 438

KTVVCTPEKVTKRMYDSFWR@FKHSEKVHVNLLLIEARM@AELLYALRAITRYMT 492

Pedant information for DKFZphamy2_l2lm2, frame 3

Report for DKFZphamy2_121m2.3

15	CLENC CMWD CpID	стна	480 54493.92 5.57					•	
20	EBLO	OLI lated -T2 ni	TREMBL:AF(PA36-T2 nuclear uclear protein PR000490	r pro	tein"; H	lomo sap	iens p53	regula	53 ated
25	EKM3		A11_A1pha LOW_COMPLEXITY		3.75 %		*	·	
30	SEQ SEG PRD		DSECRAELKDYLRFAF						
.30	SEQ SEG PRD		GLEALMSSGRVDNL						• • • • • •
35	SEQ SEG PRD		CSYLVGSHMAEFL@To		· · · · · · · ·	• • • • • • •			
40	SEQ SEG PRD		SEHTWSLAELIQALVI ncchhhhhhhhhhhhh				x x x	xxxxxx	«xxxxx
45	SEQ SEG PRD	×ו••	PLNNSGGFESARDVEA						
50°	SEQ SEG PRD		CCCCCCEEEECCCCC						
<i>3</i> 0	SEQ SEG PRD		KF@AAYSLTYNTIAM					• • • • • •	• • • • • •
55	SEQ SEG PRD		KTVACYPEKTTRRMY						

(No Prosite data available for DKFZphamy2_121m2-3)

(No Pfam data available for DKFZphamy2_121m2.3)

5

5 group: transmembrane protein

DKFZphamy2_121o17 encodes a novel 212 amino acid protein without similarity to known proteins.

- The novel protein contains I transmembrane region.
 No informative BLAST results: No predictive prosite, pfam or SCOP motife.
- The new protein can find application in studying the expression profile of amygdala-specific genes and as a new marker for amygdala cells.

unknown protein

20

Pedant: TRANSMEMBRANE 1

Sequenced by DKFZ

25 Locus: /map="186.6 cR from top of Chr22 linkage group"

Insert length: 2690 bp
Poly A stretch at pos. 2666, polyadenylation signal at pos. 2634

30 L TGCTGGGAÁA AGTGACTGCG ATTCTGAAGA ACCGCTGCCT TGCAAGGTCA 51 AGGACATTCA GTGGTTGCTG GGGTCCGCAG ACTACTGCCA CCCACTCACC DD ATCAACTCTG TTAGCCCAAT TGCCCTGCTG AACAACTGCC TGAATACAGG
151 CTTTAGGTTC CCCTGGACTC CAGCCAAGGC TGTTCAGGTG GGACCATGGT
201 GCTCTTTAAG CGTGATCGGA GGGAAGACAC ACAGCAGGGC CACCATTCCA
251 TGAATGGGAG GTGTACAGAT CACTTTCTCT TTGTGCTCAG TTCTCTTCTG
3D1 TCTCCAGCAG CTATATTGGT AAGACTAGTA CCTGCCAGGG AGAGGTGCCC 35 351 CCAAGTGAAG GGGTACAGTG GCACCTGGGA AAAGGCACCT GGAAGGTTTC 401 CATGTGGCCC AGCCCAGCAT GGAAGCAGGG TGGGAACTCT GCTGTGTCGC 451 CAGCCCTCAC TCTACTCAAG TGGCTTTTTG AGAGCCCTGC CATGTCTGTG
501 TCAGGCCTGT GCTGCTTCAC ACCCTACAGC TGCCTGGGAA AGGCCGGCCA 40 551 CGCTCCCTGT CCACACACTC CCTGTCCACA CACTCCCTGT CCACAACTGC LD1 AGCCGGGCCC TCTGCCTATG GGCACCCAAT CCAAGCAGCT GCTCCACCTT LSI TGTTTGGCAT GGTGATTTGT GTTTTTTCTC TTGGTGCTTA TGTGTGTGGG 701 CTTGGGACGA GTGCTGGTAT GCACTTAGGA CCTTCTTGAT AGCTCCCTGC 45 751 ACTTTGGAAC ACGGAGCAGA TGAGAGAGGG TCAGGGGCTT GCCCTCCACC BD1 TTGGACTTGG AAGAAGCCCA CATTGGAGAG GTGAGGACCC CATGGTGGCT BS1 CTAGTGGAAG ATACGTTAGT CTCCAGCTAA GGAGGATGAG GCGCAGCCCC PDL AGAGGGAGAC CTCAGTGATA GGGGATCAGG CTACGAAAGT GGGGGAAGGG 951 AGATGCTTTG TACATATTTT GGGGTTATAA TTTCTCTAAA TTTTAGGAGA 50 1001 ACGGGTATTG ATTGATAAAA GGGACAGGCA GTAGTGTTCA ACAGTGCATG 1051 TGAAGGAAAG TTCTGTTTTC CATGGTTTTG ACATTCTTTG GACTGTATTG LIDL TGACTGCTGT CTGGTCCACA TGGTACCCTT TTGGTAAGTA GGCTTCAGTG
LL5L CATACCAGGG TATCACTGGA GATGGGAGTT AGTGAAGGGG TGACTCCCTG
L2DL GCCTAGTATA GTGTGACCCT GGGACAACTT AATGTCCTAA AGCATTTTGG
L25L TGACTTCTAG GGAATAGCAA AGACCTATTT CATTGTCCCC AGGTAAGTAT
L3DL GTGATGAGCA ATGAGGAGGA GTGGAAAACA AAACCCAGAA AGTGCGGCAG
L35L GACCAGCCTG ACGCACACGC TCCTGTTGTC ATGGCAGACA GCCGCCTTGG *5*5.

WO 01/98454



	1401	GTGGGCACCA	CCCTGGCACT	TCCACCCTCT	AGGGGAGTGA	AGGGACATGG
	1451	CTGAGCTGGG		GGTTGACTTA	GGGAACAAGC	CCTGGGATTG
	1501	GACAAAAGGG	CCCATGCTGC		TGGGGGCAGA	
	1551	GAAGAGGGAA				GCTCTGGGTG
5			GAGATCCTAA	TGGAGGCGCC	TCCATCTGCA	ACCACAGTTG
3	7607	TAAGGCTCAT	GGCACCTCTG	CTTGGAAAGC	ACTGGTTTAG	GGACTTAGAG
	1651	AGGTAGGCAC	AAGGTGGGTC	TCCTGGGTAA	GGGAAGCAAG	AGCAGACTGT
	1701	TGGGCCAACA		CCCAGAGTAG	GGGAGAAGGT	TGGGGTGTAG
	1751	GGCCTTCCAC	GTGGAACAGA	CAGCCCCTGT	GTCTCTGTCT	CTTGGGGACC
	7907	TGAGTTTGGG	TGGGGTGGCA	GTTGGCACAG	CGCAGATGCG	GTAGAGATGG
10	1851	GAGGAAACCC	AGCTCCTCAC	TTCCGTGTGC	CTCATGCCTT	TGCATACACA
	1901	AGCACCAAAC	CTACTAGGTC	TTCTCATTAC	CCATGTAAAC	CACATGTTAG
	1951	ATAAATTTTT	GCAAGTAGAG	GAAAGAAGGA	AATAAAACAT	CACATTTTGG
	5007	TGTCTCTCAG	GCTTTCCCCC	CCAACTATGG	TTTCTTTGCT	TTTTGTTTTA
	2051	ACATAGTTTT	GTTGCTGTCT	TCTGTAATGA	TACAGTTTTG	TGCAGCTGTT
15	5707	TTCACTTAGC	ATATCGTGGG	CATCTCCCCT	TATGATTACT	AAATATTTTA
	2151	TTTTGGAGTG	GCTGTGTACT	CTCCCATTGA	CTAGATGGAC	CATTGTGCCA
	5507	GTTGCCAATC	ACTAATGCTG	TTACTAACTT	TTCAGTTATA	AATTGATGAA
	2251	TATCTTTGTG	CACAGGCTGT	TTCCCAATGT	CAAGTTATTA	GGGTAGACTC
	5307	CAGGAGGTGG	GATTCTTCAA	CTAAAGAATA	TGAAAACCTT	TGAGGCTTTT
20	2351	ACTACATATT	GACAAAATGG	TTTCCGGAAA	TATTTGTATC	CCCTTACACT
	2401	GCCACCAGCA	AGGATAAACA	TGTCCATCTT	GCCCGTATTG	GGAATTATCA
	2451	TCTGGCTAAA	TATTTGCTAA	TTTGATAATG	AAAAAATAGC	ATCGTGTTTC
	2501	AGTTGGCATT	TCACTGACTT	CTAGCACGGT	TGAACATCTT	TCATGTGGAG
	2551	CGATTGTATT	TCCTCCTTTG	TGGATTGTCA	GTGTCCTTTG	CTCTATCTTC
25	5207	TGGGGTCAGA	TAAATTTGTA	TGAGCTCGGT	ATATATTAAA	
	2651	TEGTETETET	CAAAAAAAA		AAAAAAAAA	DATATIAACC
		. 55, 5, 6, 6,	CHANANANA		ODDARABARA	

BLAST Results

30

Entry HS1033E15 from database EMBL:
Human DNA sequence from clone 1033E15 on chromosome 22q13.1-13.2.
Contains part of a novel gene, ESTs and a GSS.

35 Score = 5919, P = 5.le-262, identities = 1187/1195

Entry HSN128A12 from database EMBL: Human DNA sequence from cosmid N128A12 on chromosome 22q12-qter contains ESTs, CpG island.

40 Score = 5038, P = 0.0e+00, identities = 1014/1019

Entry HSL9034L from database EMBL: human STS WI-14034.

Score = 1800_7 P = $1.4e-76_7$ identities = 392/417

45

Medline entries

50

No Medline entry

55

Peptide information for frame ${\tt L}$

ORF from 196 bp to 831 bp; peptide length: 212

PCT/IB01/02050 Category: putative protein Classification: no clue I MVLFKRDRRE DTQQGHHSMN GRCTDHFLFV LSSLLSPAAI LVRLVPARER 51 CPQVKGYSGT WEKAPGRFPC GPARHGSRVG TLLCRQPSLY SSGFLRALPC 101 LCQACAASHP TAAWERPATL PVHTLPVHTL PVHNCSRALC LWAPNPSSCS 5 J51 TFVWHGDLCF FSWCLCVWAW DECWYALRTF LIAPCTLEHG ADERGSGACP 201 PPWTWKKPTL ER 10 BLASTP hits No BLASTP hits available 15 Alert BLASTP hits for DKFZphamy2_121017, frame 1 No Alert BLASTP hits found 20 Pedant information for DKFZphamy2_121o17, frame 1 Report for DKFZphamy2_121o17.1 25 **ELENGTHI** 575 23727.55 EMMI **Elgl** 8.73 EKWI TRANSMEMBRANE 30 SEQ MVLFKRDRREDTQQGHHSMNGRCTDHFLFVLSSLLSPAAILVRLVPARERCPQVKGYSGT PRD MEM 35 SEQ WEKAPGRFPCGPAQHGSRVGTLLCRQPSLYSSGFLRALPCLCQACAASHPTAAWERPATL PRD MEM 40 SEQ PVHTLPVHTLPVHCSRALCLWAPNPSSCSTFVWHGDLCFFSWCLCVWAWDECWYALRTF PRD CCCCCCCCCCCCeeeeeccccceeecccceeecccceeeccchhhhhhhe MEM SEQ LIAPCTLEHGADERGSGACPPPWTWKKPTLER 45 PRD eeecccccccccccccccccccccc MEM (No Prosite data available for DKFZphamy2_121o17.1) 50

WO 01/98454

(No Pfam data available for DKFZphamy2 121o17.1)

DKFZphamy2_12d7

5 group: signal transduction

DKFZphamy2_12d7 encodes a novel 552 amino acid protein, which is a so far unknown alternative spliced form of disks large homolog DLG2.

10

15

20

It seems to be predominantly expressed in the retinar germ cells and brain. It contains a SH3-domain and a guanylate kinase domain. These conserved regions are shared among members of the discs-large family of proteins that include human p55; a membrane protein expressed in erythrocytes; rat PSD-95/SAP90; a synapse protein expressed in brain; Drosophila dIg-A; a septate junction protein expressed in various epithelia; and human and mouse Z0-1 and canine Z0-2; two tight junction proteins. The Homologue of Drosophila; dIg-A; acts as a tumor suppressor. All members of this family may be involved in signal transduction.

The new protein can find application in modulating/blocking intracellular signal transduction pathways.

25

similarity to disks large homolog DLG2 (Homo sapiens)

alternative splicing: see DLG2 complete cds.

30 frame shift: around position 1437 one C too many

Sequenced by EMBL

Locus: /map="338.6 cR from top of Chrl7 linkage group"

35

Insert length: 4220 bp

Poly A stretch at pos. 4160, polyadenylation signal at pos. 4165

1 CCCGGCTGCG CTGGAGCCGC CCGGAGCTAG GGGCTTCCCG GGGCGCAGGA
51 GAGACGTTTC AGAGCCCTTG CCTCCTTCAC CATGCCGGTT GCCGCCACCA 40 101 ACTCTGAAAC TGCCATGCAG CAAGTCCTGG ACAACTTGGG ATCCCTCCCC
151 AGTGCCACGG GGGCTGCAGA GCTGGACCTG ATCTTCCTTC GAGGCATTAT
201 GGAAAGTCCC ATAGTAAGAT CCCTGGCCAA GGCCCATGAG AGGCTGGAGG 251 AGACGAAGCT GGAGGCCGTG AGAGACAACA ACCTGGAGCT GGTGCAGGAG
301 ATCCTGCGGG ACCTGGCGCA GCTGGCTGAG CAGAGCAGCA CAGCCGCCGA
351 GCTGGCCCAC ATCCTCCAGG AGCCCCACTT CCAGTCCCTC CTGGAGACGC 45 401 ACGACTCTGT GGCCTCAAAG ACCTATGAGA CACCACCCC CAGCCCTGGC 451 CTGGACCCTA CGTTCAGCAA CCAGCCTGTA CCTCCCGATG CTGTGCGCAT 501 GGTGGGCATC CGCAAGACAG CCGGAGAACA TCTGGGTGTA ACGTTCCGCG 50 551 TGGAGGGCGG CGAGCTGGTG ATCGCGCGCA TTCTGCATGG GGGCATGGTG
LD1 GCTCAGCAAG GCCTGCTGCA TGTGGGTGAC ATCATCAAGG AGGTGAACGG **L51 GCAGCCAGTG GGCAGTGACC CCCGCGCACT GCAGGAGCTC CTGCGCAATG** 701 CCAGTGGCAG TGTCATCCTC AAGATCCTGC CCAGCTACCA GGAGCCCCAT 751 CTGCCCCGCC AGGTATTTGT GAAATGTCAC TTTGACTATG ACCCGGCCCG 55 BOD AGACAGCCTC ATCCCCTGCA AGGAAGCAGG CCTGCGCTTC AACGCCGGGG B51 ACTTGCTCCA GATCGTAAAC CAGGATGATG CCAACTGGTG GCAGGCATGC 901 CATGTCGAAG GGGGCAGTGC TGGGCTCATT CCCAGCCAGC TGCTGGAGGA

PCT/IB01/02050

	951	GAAGCGGAAA	GCATTTGTCA	ACACCCACCT	GGAGCTGACA	·CCAAACTCAC
	1001				AAAAGAAGCG	_
	1051	TTGACCACCA				AATGATGTAT
	1101	GGAGGTGGCC		GTTTGACCGT	CATGAGCTGC	TCATTTATGA
5					GAAAACCCTG	GTACTGATTG
3	1151	GGGCTCAGGG			AGAACAAGCT	CATCATGTGG
	7507	GATCCAGATC		·CACGGTGCCC	TACACCTCCC	GGCGGCCGAA
	1527	AGACTCAGAG		AGGGTTACAG	CTTTGTGTCC	CGTGGGGAGA
	7307	TGGAGGCTGA	CGTCCGTGCT	GGGCGCTACC	TGGAGCATGG	CGAATACGAG
	1351	GGCAACCTGT	ATGGCACACG	TATTGACTCC	ATCCGGGGCG	TGGTCGCTGC
10	1401	TGGGAAGGTG	TGCGTGCTGG	ATGTCAACCC	CCAGGCCGGT	GAAGGTGCTA
	1451	CGAACGGCCG	AGTTTGTCCC	TTACGTGGTG	TTCATCGAGG	CCCCAGACTT
	1501	CGAGACCCTG	CGGGCCATGA	ACAGGGCTGC		GGAATATCCA
	1551	CCAAGCAGCT	CACGGAGGCG		GGACAGTGGA	GGAGAGCAGC
	7207	CGCATCCAGC			GACCTCTGCC	TGGTCAATAG
15	1651	CAACCTGGAG				
~-	1701	GGACAGAGCC			GACAGCCATG	GAGAAGCTAC
•	1751	CCTGGTCCTT		CCTGTCAGCT	GGGTGTACTG	AGCCTGTTCA
			GGCTCACTCT	GTGTTGAAAC	CCAGAACCTG	AATCCATCCC
	1801	CCTCCTGACC		TGCCACAATC	CTTAGCCCCC	ATATCTGGCT
20	1851	GTCCTTGGGT	AACAGCTCCC	AGCAGGCCCT	AAGTCTGGCT	TCAGCACAGA
20	1901	GGCGTGCACT	GCCAGGGAGG	TGGGCATTCA		TGTGCCCAGG
	1951	TGCTGCCCAC	TCCTGATGCC	CATTGGTCAC	CAGATATCTC	TGAGGGCCAA
	5007	GCTATGCCCA	GGAATGTGTC	AGAGTCACCT	CCATAATGGT	CAGTACAGAG
	2051	AAGAGAAAAG	CTGCTTTGGG	ACCACATGGT	CAGTAGGCAC	ACTGCCCCTG
	5707	CCACCCCTCC	CCAGTCACCA	GTTCTCCTCT	GGACTGGCCA	CACCCACCC
25	2151	ATTCCTGGAC	TCCTCCCACC	TCTCACCCCT	GTGTCGGAGG	AACAGGCCTT
	5501	GGGCTGTTTC	CGTGTGACCA	GGGGAATGTG	TGGCCCGCTG	GCAGCCAGGC
	2251	AGGCCCGGGT	GGTGGTGCCA	GCCTGGTGCC	ATCTTGAAGG	CTGGAGGAGT
	5307	CAGAGTGAGA	GCCAGTGGCC	ACAGCTGCAG	AGCACTGCAG	CTCCCAGCTC
	2351	CTTTGGAAAG	GGACAGGGTC	GCAGGGCAGA	TGCTGCTCGG	TCCTTCCCTC
30	2401	ATCCACAGCT	TCTCACTGCC	GAAGTTTCTC		CAATGTGTCC
	2451	TGACAGGTCA	GCCCTGCTCC	CCACAGGGCC	AGGCTGGCAG	GGGCCATTGG
	2501	GCTCAGCCCA	GGTAGGGGCA	GGATGGAGGG	CTGAGCCCTG	
	2551	CTGTTACCAA	CTGAAGAGCC	CCAAGCTCTC	CATGGCCCAC	
	5207	GGTCTGAGCT	CTATGTCCTT	GACCTTGGTC	CATTTGGTTT	AGCAGGCACA
35	2651	CAGGTCCAGG	TAGCCCACTT	GCATCAGGGC		TCTGTCTAGC
.55	2707	GGAGGAGTGC	AGAGGGGACC		TECTEGETTE	GAGGGGCTAA
	2751	CTCCAGGAGG	TTCCTCACAC	TTGGGAGCCT	GGGCTTGAAG	GACAGTTGCC
	5807	TCTGTACAAC		ACAACTCCAG	AGGCGCCATT	TACACTGTAG
	2851		CTGTGGTTCC	ACGTGCATGT	TCGGCACCTG	TCTGTGCCTC
40		TGGCACCAGG	TTGTGTGTGT	GTGCGTGTGC	ACGTGCGTGT	GTGTGTGTGT
40	2907	GTGTCAGGTT	TAGTTTGGGG	AGGAAGCAAA	GGGTTTTGTT	TTGGAGGTCA
	2951	CTCTTTGGGG	CCCCTTTCTG	GGGGTTCCCC	ATCAGCCCTC	ATTTCTTATA
	3007	ATACCCTGAT	CCCAGACTCC	AAAGCCCTGG	TCCTTTCCTG	ATGTCTCCTC
			TTGTCCCCCT	ACCCTAAATG	CCCCCTGCC	ATAACTTGGG
		GAGGGCAGTT	TTGTAAAATA	GGAGACTCCC	TTTAAGAAAG	AATGCTGTCC
45		TAGATGTACT	TGGGCATCTC	ATCCTTCATT	ATTCTCTGCA	TTCCTTCCGG
		GGGGAGCCTG	TCCTCAGAGG	GGACAACCTG	TGACACCCTG	AGTCCAAACC
	3251	CTTGTGCCTC	CCAGTTCTTC	CAAGTGTCTA		GCTGCAGCGT
		CAGCCAAAGC	TGGCCCCTGA	ACCACTGTGT	GCCCATTTCC	TAGGGAAGGG
		GAAGGAGAAT	AAACAGAATA	TTTATTACAA	ATGTTAGAAT	ATATTTCTTA
50		TACTAGGAAT				GGGGTGGAAA
-		GGCCAGGCCT		CGTTGGTGTG		
	3501	TCATTCTCCT	GCTCCTCTTT			TACTACACAC
		CAGCTCTGCC		. –		CATCCTGATT
			TTGCATCACC			AGGAAATGGG
55	21 L2 TP G T	TCTCTCCCCC	GCTGACCTGA	DUCTATAGGG	TCACTTGCCA	TTTCCTACCT
55	3651	1C1C166666	ATTTGAGGGT	AGAGGCAGGG		TGTTGCAGTT
	3701	GCTTCTGCCC	CCTTGATCCA	AATGACCATC		GAGATGGGTT
	3751	GGGTACCTGG	CCTTCATGGC			TCAAGGGGCA
	3907	GGCCTGGGGC	CCTTCCCTCC	TGTCTCTTCT	CGGTCTTTCC	TCTCTGAGCA

WO 01/98454

PCT/IB01/02050

3851 GCCTCCTACC TCCCCTGCCT GAGCCCTCAC TCCACAGCCC TCCCAGGTAC
3901 CTAGCAGAGG CTGTCAGTCC TTGGCTCACC TGGAACAGGG CTGGGGCTGG
3951 GTTGGAACAG GTGTGTGCCC CCACCACAGC TCTATGACTC TGTTCTCCCT
4001 CCCTGCCATT GTGGACTCTT GTATTTGAGG GACCTCAAGA GAGTGAGGAC
5 4051 CCTACCATCC ACTGTCCATA TTCAGTCCCA GCCCCAGTGC GCTTCCTCTG
4101 TTCCCTCCCT CAGCCATCCA ATTCTTGAGT TTTCTCACTG ATTGGTTTTC
4151 TTTCTTTTTC CTTGGATTAA ATGTGAAAGC AAAGAAAAAA AAAAAAAAA

10

BLAST Results

No BLAST result

15

Medline entries

20 96070428:

Mazoyer S, Gayther SA, Nagai MA, Smith SA, Dunning A, van Rensburg EJ,

Albertsen Ha White Ra

Ponder BA: A gene (DLG2) located at 17q12-q21 encodes a new

25 homologue

of the Drosophila tumor suppressor dIg-A. Genomics 1995 Jul

1:28(1):25-31

30

Peptide information for frame 1

35

ORF from 82 bp to 1437 bp; peptide length: 452 Category: strong similarity to known protein Classification: Cell signaling/communication Prosite motifs: GUANYLATE_KINASE_1 (385-402)

40

	1	MPVAATNSET	AMQQVLDNLG	SLPSATGAAE	LDLIFLRGIM	ESPIVRSLAK
			EAVRDNNLEL			
			ASKTYETPPP			
45			ELVIARILHG			
			VILKILPSYQ			
•			IVNQDDANWW			
			GZLZGKKKKR			
			VGRRSLKNKL			
50			VRAGRYLEHG			
	451					

55

BLASTP hits

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_12d7, frame 1

No Alert BLASTP hits found

5

Peptide information for frame 2

ORF from 1439 bp to 1738 bp; peptide length: 100
10 Category: strong similarity to known protein
Classification: Cell signaling/communication
Prosite motifs: LEUCINE_ZIPPER (66-87)

15 L VKVLRTAEFV PYVVFIEAPD FETLRAMNRA ALESGISTKQ LTEADLRRTV 51 EESSRIQRGY GHYFDLCLVN SNLERTFREL QTAMEKLRTE PQWVPVSWVY

20

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_12d7, frame 2

BLASTP hits

No Alert BLASTP hits found

Pedant information for DKFZphamy2_12d7, frame 1

30

EPIRKUJ

EPIRKUJ

EPIRKWI

[PIRKW]

[PIRKU]

55

25

Report for DKFZphamy2_12d7.1

ELENGTHD 516 35 EMMI 56458.36 [[q] P-57 **EHOMOLI** PIR:A57653 disks large homolog DLG2 - human D.D EFUNCATI 01.03.99 other nucleotide-metabolism activities cerevisiae, YDR454cl 7e-15 **EFUNCATI** f nucleotide metabolism and transport EH. influenzae. HI17431 3e-07 **EBLOCKSI** PRODB34F EBF0CK21 BL00856C BLOOA56B Guanylate kinase proteins **EBLOCK21 EBLOCKSI** BLOOB56A Guanylate kinase proteins [[SCOP]] dlgky___ 3.29.1.1.1 Guanylate kinase yeast (Saccharomyce &e-45 EZC0P3 dlkwab_ 2.26.1.1.2 Cask/Lin-2 [Human (Homo 4e-34 sapiens) 50 EECI 2.7.4.8 Guanylate kinase 8e-17 **EPIRKUJ** blocked amino end 8e-17 **EPIRKWI** phosphotransferase &e-17 **EPIRKU** monomer &e-17

-82-

duplication 5e-29

P-loop Be-17

signal transduction 3e-24

alternative splicing 5e-29

acetylated amino end le-16



PD0CB0670

	WO 01/98454				PCT/I	B01/02050	
	EPIRKWI EPIRKWI EPIRKWI	magnesiu ATP 8e-l	?	<u>-</u> 74			
.	ESUPFAMD ESUPFAMD ESUPFAMD	SH3 homology discs-large t unassigned Se	umor suppre	essor 3e-24 /r-specific	protein	kinases	5e-
10	EMATQUES EMATQUES EMATQUES EMATQUES	protein kinas GLGF domain h guanylate kin guanylate kin	omology 9e- ase 8e-17	-74			
15	EPFAMI EKWI	GUANYLATE_KIN					
15	EKWI	עב		٠.			
20	SEQ MPVA	ATNSETAMQQVLDN					
20	laky-	DNNLELVÆEILRDL	AGLAEGSSTAA	\ELAHIL@EPH	FQSLLETHD	SVASKTYE	тррр
25	SEQ SPGLI	PTFSN@PVPPDAV	RMVGIRKTAGE	THLGVTFRVEG	GELVIARIL	HGGMVAQQ	GLLH
		• • • • • • • • • • • • • • • • • • • •					
30	lgky-	CKEVNGQPVGSDPR			•		
35	SE@ DSLIF	PCKEAGLRFNAGDL				*	
40	lakv-	ISGTLCGSLSGKKK					
	SEQ VGRRS	LKNKLIMUDPDRY	GTTVPYTSRRF	KDSEREGQGY	SFVSRGEME	ADVRAGRY	LEHG
45		ІНННННННТТТТЕ					
	lgky-	EEEEEHHHHHHHH					
50	SEQ QGCAG	EUNIHQAAHGGGP					
55		Pro	osite for D	KFZphamy2_	12d7-1		
	P20085L	385->403	GUANYLATE	_KINASE_l	PDO	C BO6 70	

GUANYLATE_KINASE_1

Pfam for DKFZphamv2 12d7.1

```
5
    HMM_NAME Src homology domain 3
    HMM
    *pyVIALYDYqAqd.....pDELSFkEGDIIiIIEdsDD-WWrgRnnn
10
                      +V+ +DY++ +
                                      + LF GD ++I+++D+ WW +
                955
    VFVKCHFDYDPARDSLIPCKEAGLRFNAGDLLQIVNQDDANWWQACHVE 276
                    TNGQEGWIPSNYVEP:*
15
                    ++ G+IPS +E+
    Query
                277 GG-SAGLIPSQLLEEK
                                     297
20
              Pedant information for DKFZphamy2_12d7, frame 2
                      Report for DKFZphamy2_12d7.2
25
   ELENGTHD
           175
   EMMI
            19721.90
   [[q]
            9.69
   EHOMOLI
                PIR:A57653 disks large homolog DLG2 - human 7e-53
30
   EPIRKWI
                membrane protein le-13
   ESUPFAMI
            SH3 homology le-13
   CSUPFAMI
            GLGF domain homology le-13
   ESUPFAMI
            guanylate kinase homology le-13
   EPROSITED LEUCINE_ZIPPER 1
.35
   [KW]
            Alpha_Beta
   SEQ
       MAPRCPTPPGGRKTQSGKVRVTALCPVGRWRLTSVLGATWSMANTRATCMAHVLTPSGAW
   PRD
       40
   SEQ
       SLLGRCACUMSTPRPVKVLRTAEFVPYVVFIEAPDFETLRAMNRAALESGISTKOLTFAD
   PRD
       SEQ
       LRRTVEESSRIQRGYGHYFDLCLVNSNLERTFRELQTAMEKLRTEPQWVPVSWVY
45
   PRD
       Prosite for DKFZphamy2_12d7.2
50
   PS00029
               141->163
                        LEUCINE_ZIPPER
                                              PD000029
   (No Pfam data available for DKFZphamy2_12d7.2)
55
```

WO 01/98454 DKFZphamy2_12g7

5 group: amygdala derived

DKFZphamy2_12g7 encodes a novel 254 amino acid protein without similarity to known proteins.

No informative BLAST results: No predictive prosite: pfam or SCOP motife.

The new protein can find application in studying the expression profile of amygdala-specific genes.

putative protein

Sequenced by EMBL

20 Locus: unknown

15

25

Insert length: 1257 bp

No poly A stretch found, no polyadenylation signal found

L CTCCAAGACT TCCTTGCTGT GAGGCTCGTG TGGACCCCAG AGCATGCACA 51 GGCTGTTTAC TCCACAGAGT GGCTTTGAGA ATCAGATGAG ACTGTGCTGG 101 CGAAGGCCCT GTGGGAATGA GGAACGCTGT AGTGTTTGCT GGTCCCTGTT 30 151 TCTGCCCCA GGAAAGCAGC TGTGTGAGGA GGAGCGCCGG GCCATGCAGG 201 CTGCCCTGGA CTCCGTCGTC TGCCACACGC CCCTCAACAA CCTTGGCTTT 251 TCCCGGAAGG GCAGCGCGCT CACCTTCAGT GTGGCCTTCC AGGCTCTGAG
301 GACGGGGCTC TTCGAGCTAA GCCAGCACAT GAAACTGAAG CTGCAGTTCA 35% CCGCCAGCGT GTCCCACCCT CCACCCGAGG CCCGGCCCCT CTCCCGCAAG 4D1 AGCAGCCCCA GAAGCCCTGC TGTCCGGGAC TTGGTGGAGA GGCATCAGGC 35 45% TAGCCTGGGC CGCTCCCAGT CCTTCTCCCA CCAGCAGCCT TCCCGAAGCC 501 ACCTCATGAG GTCGGGCAGT GTGATGGAGC GCAGAGCATC ACGCCCCCTG
551 TGGCCTCTCC TGTTGGCCGC CCCCTCTACC TGCCCCCGGA CAAGGCTGTG
601 TTGTCTCTGG ACAAGATTGC CAAGCGCGAG TGCAAGGTCC TGGTGGTGGA
651 ACCCGTCAAG TAGCACCGTG CCAGCTCTGT TCCCTCTTAC ACTCCAGAGA
701 CCCAACGCCC CCAGAGGGTA TCCTTGCTCC CGGGCTGTGC CTCCCTGGG 40 751 ATGCCTCCCA GACGGGGGTG AAGAGGCCTG GCAGAGCTGC CTGTCTTGTG BOL TCTGCTGATG AGGGATGGGG GAAGAAGCTG TGAAGTGGGC GGGCATGGCT BSL GGGACTAAGC CACCAGTATT CCCCGACGTT CCTGTGGGGG GGGCTGGCCC 901 ACCCCTAGGC CAGGGCAAGG GTTCCCAGAG CTCCCTTGTC CCCGGCCCTT 45 951 TACCCTGGTT CTGAGTTTAC AAAGTCTCTT CCTCATTCCC GTTGAGTTCT BODD TICCCACCIC TGACATICCC TCCCTCCCTC CCGCAGGCTG AGATTAGAGG 1051 GTGGTGATGG CTAAGGGCCC CTGACAGTGA CCTTCCTGTC TCAGGGGTTG 1101 GGGACAGGGC CAGGTAGCCT CCTGCCCCTT ATGTTTACGT TTGCAGCCTG 1151 AAGCACTTTA ATTTTTTTT TTTTTGGTCT GTCCCTGTAA CTAATTTTCC 50 1201 AACTATTGCT TCCAACTGAA ATAAGACTAT TAAATGCCTG TTCAGAGGGA 1251 AAAAAA

55 BLAST Results

No BLAST result

SEQ

ZEG

PRD



Medline entries

5 No Medline entry 10 Peptide information for frame 2 ORF from 44 bp to 805 bp; peptide length: 254 Category: putative protein Classification: no clue 15 1 MHRLFTPQSG FENQMRLCWR RPCGNEERCS VCWSLFLPPG KQLCEEERRA 51 MQAALDSVVC HTPLNNLGFS RKGSALTFSV AFQALRTGLF ELSQHMKLKL 101 QFTASVSHPP PEARPLSRKS SPRSPAVRDL VERHQASLGR SQSFSHQQPS 151 RSHLMRSGSV MERRASRPLW PLLLAAPSTC PRTRLCCLWT RLPSASARSW 20 201 WWNPSSSTVP ALFPLTLQRP NAPRGYPCSR AVPPLGCLPD GGEEAWQSCL 527 ZCAC 25 BLASTP hits No BLASTP hits available 30 Alert BLASTP hits for DKFZphamy2_12g7, frame 2 No Alert BLASTP hits found Pedant information for DKFZphamy2_12g7, frame 2 35 Report for DKFZphamy2_12g7.2 ELENGTHD 254 40 28479.91 EMWD 10.00 EBLOCKSI BL01013C Oxysterol-binding protein family proteins Alpha_Beta 4.72 % 45 LOW_COMPLEXITY ZEQ MHRLFTPQSGFENQMRLCWRRPCGNEERCSVCWSLFLPPGKQLCEEERRAMQAALDSVVC SEG 50 PRD HTPLNNLGFSRKGSALTFSVAFQALRTGLFELSQHMKLKLQFTASVSHPPPEARPLSRKS SEQxxxxxx SEG PRD 55

SPRSPAVRDLVERHQASLGRSQSFSHQQPSRSHLMRSGSVMERRASRPLWPLLLAAPSTC

xxxx......

	SEQ	PRTRLCCLWTRLPSASARSWWWDSSSTVPALFPLTLQRPNAPRGYPCSRAVPPLGCLPD
	SEG	• • • • • • • • • • • • • • • • • • • •
	PRD	CCCCEEEECCCCCCCCEEECCCCCCCCCCCCCCCCCCCC
5	SEQ SEG PRD	GGEEAWQSCLSCVC
10	(No	Prosite data available for DKFZphamy2_12g7.2)
. •	(No	Pfam data available for DKFZphamy2_12g7-2)

DKFZphamy2_12i1

5 group: amygdala derived

DKFZphamy2_12il encodes a novel 283 amino acid protein with weak similarity to F41E6.3 of Caenorhabditis elegans.

No informative BLAST results; No predictive prosite, pfam or SCOP motife.

The new protein can find application in studying the expression profile of amygdala-specific genes.

putative protein

Sequenced by EMBL

Locus: /map="3"

15

20

25

Insert length: 2528 bp

Poly A stretch at pos. 2515, polyadenylation signal at pos. 2491

L ATATAGTTGG ATCAAACAAA AACAACACAA TTTGTCCCGA TAATTATCAA 51 ACAGCACAGC TACTTGCCTT AATTTTAGAG TTACTCACAT TTTGTGTGGA LOL ACATCACACA TATCACATAA AAAACTATAT TATGAACAAG GACTTGCTAA 30 151 GAAGAGTCTT GGTCTTGATG AATTCAAAGC ACACTTTTCT GGCCTTGTGT 201 GCCCTTCGCT TTATGAGGCG GATAATTGGA CTTAAAGATG AATTTTATAA 251 TCGTTACATC ACCAAGGGAA ATCTTTTTGA GCCAGTTATA AATGCACTTC BOL TGGATAATGG AACTCGGTAT AATCTGTTGA ATTCAGCTGT TATTGAGTTG 351 TTTGAATTTA TAAGAGTGGA AGATATCAAG TCTCTTACTG CCCATATAGT 401 TGAAAACTTT TATAAAGCAC TTGAATCGAT TGAATATGTT CAGACATTCA 35 451 AAGGATTGAA GACTAAATAT GAGCAAGAAA AAGACAGACA AAATCAGAAA 501 CTGAACAGTG TACCATCTAT ATTGCGTAGT AACAGATTTC GCAGAGATGC 551 AAAAGCCTTG GAAGAGGATG AAGAAATGTG GTTTAATGAA GATGAAGAAG LOI AGGAAGGAAA AGCAGTTGTG GCACCAGTGG AAAAACCTAA GCCAGAAGAT 40 LSI GATTTTCCAG ATAATTATGA AAAGTTTATG GAGACTAAAA AAGCAAAAGA 701 AAGTGAAGAC AAGGAAAACC TTCCCAAAAG GACATCTCCT GGTGGCTTCA 751 AATTTACTTT CTCCCACTCT GCCAGTGCTG CTAATGGAAC AAACAGTAAA BOL TCTGTAGTGG CTCAGATACC ACCAGCAACT TCTAATGGAT CCTCTTCCAA B51 AACCACAAAC TTGCCTACGT CAGTAACAGC CACCAAGGGA AGTTTGGTTG 45 901 GCTTAGTGGA TTATCCAGAT GATGAAGAG AAGATGAAGA AGAAGAATCG 951 TCCCCCAGGA AAAGACCTCG TCTTGGCTCA TAAAATATTT ATTAGGGGAC DOD CCTCAACATG TGGTCTTACA ATGCTGCAAC TGTTCAGTGA GCTGAAAATC 1051 TGAATCAGAA AGCTTTCTCA ATTGAACTTA TAAAATATAC AAGGAGTAGC LLOL AAAAGACAGT ATATCAGCTA AGAGAGTTTA GTTCTAATAA AAATCAGGCT LISL TCCCAGGAAC TTGATTGCTT GCTAGTAATT AAGGGGTTTG CCTTTTAGGC
L201 TGTCAAAACA AACATTAGTA ACCAGAACCT GGGAGGTAGC TTCTCAGCAA
L251 GGAAAAGTCA CAGGTTTGGG GACGGTTTAG GGGAGGGGAA AAGGTTGATA
L301 TAATAATGCA GGGTTGCTCC TCGGGGTGTC GATCTAGAAA CAATTTTACA
L351 GAACTTCAGT TGTAAACTCA ATAACATTAC TTGTATAATG GTGCTGGCCA
L401 TGTTGTTGTT TTAATCAGTT GCCTCTTTTT AAAAGAAATT TTTATGGAAA
L451 ACACATTCAA CTATCAATTAA AAAAATGAAG TTAAGCTGTT GCCAGGATTT 50 55 1451 ACACATTCAA CTATCATTAA AAAAATGAAG TTAAGCTGTT GGGACCATTT 1501 CTTTAAGATT TAACAAAAGT TCAGCCTTTT AGGTAGTTGA AGGGAAGTAC 1551 ACCCCGTATT CAGCACATGT TGAGTTTTCT ACACCAGGAA TTTTCAATAT

WO 01/98454

PCT/IB01/02050

	1601	GTATATTGAT	GAAAACAAGC	TCAATTCAAA	CTGGACAGTT	TTAAGATAAT
	1651	GTTAAAATCA	GCACTTTTAG	AGACAACGAA	GGCCAAGAAT	CAGTACAGTA
	1701	GTATTCCAAA	ATGATTTTCT	CTAGAAATTT	GAAAGTAGAT	CGAACAGAAT
	1751	GTTGTCAACC	GCCTACCAGT	ACAATCTTTT	GTGGAAGATA	CTTTGAAATC
5	1801	ACTTTCTACT	TTGTTAGTAA	AGTTCTGTCT	TTCCAGAGCT	GCAAGTTTTA
	1851	AAGTGTTACT	TATACAGACC	AACCAAGAAT	AGTGCTGAAT	TAAGTGGCAT
	1901	TTAGTATCTA	GAAGCCATTT	TGATCCAAGA	AGCTACTTAA	GTGTCAAAGT
	1951	CAGCATGCAG	CACATGTAGC	TTTTCTGTAA	ACAAGGGTGT	GATATGAAAG
	5007	CTGCTTTTTT	AAGAAGAGTA	AAAGCACATT	CCATATACGT	AAGTGAATTT
10	5027	TAAAAATAAA	TTGAGGCAAA	CAGTTAAGTT	TTATTTTTAG	AGCAACAAGT
	5707	TAACTGTAAA	TATTTTAATG	TTAGTTTGCT	CATCTATGAT	CTGAGATCAT
	2151	GCCGAAGTGA	GAAAAATCTC	CCCAAAATAC	AATTTAATGC	ATTĠGGAAAA
. •	5507	AAAAACTTTA	ACAGTAATTC	CAGCCACAAT	CTTTAGATCA	CCCTTGTAAT
	2251	GTGTTACGGG	TCCATTTTTC	CTGGAATCGT	TTAATCTAAA	GCAGTTTCCC
15	5307	CTGTTTTGGA	GATTTTGTAG	TTAATTTTAA	TTTTGGCTAT	TGTTTGGAAA
	2351	AGATGAGCTG	TCTGTGTAGA	TATGAAGTAT	AGTTTTTTCC	ATAAAACAGA
	2401	TGTTTATTTT		ATACCACTGT		
	2451		TGATATTAAT		AAAATTCAGG	AATTAAAATG
	- 2501	TGACCCTGTA	ATTCCAAAAA	AAAAAAA		
20						

20

BLAST Results

25 Entry AFOl6448_8 from database TREMBL:
gene: "F4lE6.3"; Caenorhabditis elegans cosmid F4lE6.
Score = 390, P = 5.0e-32, identities = 73/184, positives = 118/184,
frame +3

30

Entry HS211256 from database EMBL: human STS SHGC-15844. Score = 977, P = 5.5e-38, identities = 199/202

35

Medline entries

40 No Medline entry

Peptide information for frame 3

45

ORF from 132 bp to 980 bp; peptide length: 283 Category: putative protein Classification: no clue

50

55

I MNKDLLRRVL VLMNSKHTFL ALCALRFMRR IIGLKDEFYN RYITKGNLFE
51 PVINALLDNG TRYNLLNSAV IELFEFIRVE DIKSLTAHIV ENFYKALESI
101 EYVQTFKGLK TKYEQEKDRQ NQKLNSVPSI LRSNRFRRDA KALEEDEEMW
151 FNEDEEEEGK AVVAPVEKPK PEDDFPDNYE KFMETKKAKE SEDKENLPKR
201 TSPGGFKFTF SHSASAANGT NSKSVVAQIP PATSNGSSSK TTNLPTSVTA
251 TKGSLVGLVD YPDDEEEDEE EESSPRKRPR LGS

BLASTP hits

No BLASTP hits available 5 Alert BLASTP hits for DKFZphamy2_12il, frame 3 No Alert BLASTP hits found 10 Pedant information for DKFZphamy2_12il, frame 3 Report for DKFZphamy2_12i1.3 15 ELENGTHI 35P 37261-10 EMMI 5.60 [pI] TREMBL: AFO16448_8 gene: "F41E6.3"; Caenorhabditis **EHOMOL** 20 elegans cosmid F41Eb. le-3b EFUNCATI 01.05.04 regulation of carbohydrate utilization EScerevisiae, YNL201cJ 2e-08 EBF0CK23 BLDD357 Histone H2B proteins EBF0CK2] BP02232B 25 [BLOCK2] PRO1073C **EBFOCK2** BP03050C **EBFOCK2** BP03580F **EBFOCK2** PR00893F EKWI All_Alpha 30 EKWI LOW_COMPLEXITY 10.43 % IVGSNKNNTICPDNYQTAQLLALILELLTFCVEHHTYHIKNYIMNKDLLRRVLVLMNSKH SEQ SEG 35 PRD TFLALCALRFMRRIIGLKDEFYNRYITKGNLFEPVINALLDNGTRYNLLNSAVIELFEFI SEQ SEG PRD 40 RVEDIKSLTAHIVENFYKALESIEYVQTFKGLKTKYEQEKDRQNQKLNSVPSILRSNRFR SEQ SEG PRD 45 SEQ RDAKALEEDEEMWFNEDEEEEGKAVVAPVEKPKPEDDFPDNYEKFMETKKAKESEDKENL SEG PRD SEQ PKRTSPGGFKFTFSHSASAANGTNSKSVVAQIPPATSNGSSSKTTNLPTSVTATKGSLVG 50 SEG PRD SEQ LVDYPDDEEEDEEEESSPRKRPRLGS SEG 55 PRD eecccccchhhhhcccccccccc

(No Prosite data available for DKFZphamy2_12i1.3)

(No Pfam data available for DKFZphamy2_12i1.3)

DKFZphamy2_13g19

group: amygdala derived

DKFZphamy2_13g19 encodes a novel 281 amino acid protein without similarity to known proteins.

10 The novel protein contains a PROSITE ASP_PROTEASE motif and seem to be expressed Ubiquitously. No informative BLAST results; No predictive prosite; pfam or SCOP

The new protein can find application in studying the expression profile of amygdala-specific genes.

20 unknown protein

15

perhaps complete cds. Pedant: SIGNAL_PEPTIDE

25 Sequenced by EMBL

Locus: /chromosome="12p13.3"

Insert length: 2754 bp

30 Poly A stretch at pos. 2743, polyadenylation signal at pos. 2724

1 GCAATCTCGG GAAATTGGAG ACTGACGCGG CTGCTCCTGC ATGTTATTTA

L GCAATCTCGG GAAATTGGAG ACTGACGCGG CTGCTCCTGC ATGTTATTTA

51 TTTTTCCTCT TTCCCTCCG TGGAGACCCT CCTGTTGGAA AGAGAGCTGC

101 AGCACGGGAC AGAGACAGGC AGGAAGAAGC AGAGAGGACT CGGTGACGCC

151 CCCACCGAGC AGCCCCTGGC CCACTCCTC AGCAGGGGCC ATGAGCACCA

201 AGCAGGAGGC CAGGAGAGAT GAGGGAGAAG CCAGGACGAG GGGCAGGAG

251 GCACAGCTTC GAGACCGAGC CCACCTGAGC CAGCAGCGCC GGCTCAAACA

301 GGCCACCCAG TTCCTGCACA AGGACTCGGC CGACCTGCTC CCGCTGGACA

351 GCCTCAAGAG GCTCGGCACC TCCAAGGACT TGCAGCCGCG CAGTGTGATC

401 CAAAGACGCC TGGTGGAGGG AAACCCGAAT TGGCTTCAGG GGGAGCCTCC

451 CCGGATGCAG GACCTGATTC ATGGCCAGGA GAGCAGGAGG AAGACCAGCA

501 GGACAGAGAT TCCAGCTCTT CTGGTCAACT GCAAGTGCCA GGACCAGCTG

552 CTTAGAGTGG CCGCTTGACAC AGGCACCCAA TACAATCGGA TCTCTGCTGG

153 GGGACCTGGC CCCCTGGGGT TAGAGAAAAG GGTCCTAAAAA GCCTCAGCTG

154 GGGACCTGGC CCCCTGGGCCC CCCAACCCCAG TGGAGCAGTT GGAGCTACAG

155 GGGACCTGGC CCCTGGGCCC CCCAACCCCAG TGGAGCAGTT GGAGCTACAG

156 GGGACCTGGC CCCTGGGCCC CCCAACCCCAG TGGAGCAGTT GGAGCTACAG

157 CTGGGGCAGG AGACTTGTGGT GTGCTCGGCA CAGGTGGTG ATGCTGAGAG 35 40 45 701 CTGGGGCAGG AGACTGTGGT GTGCTCGGCA CAGGTGGTGG ATGCTGAGAG 751 TCCTGAATTC TGCCTGGGCC TGCAGACTCT GCTTTCTCTC AAGTGCTGCA BDL TCGACCTGGA GCACGGAGTG CTGCGGCTGA AAGCCCCGTT CTCAGAGCTA 50 **B51 CCCTTCCTGC CTTTGTACCA AGAGCCTGGC CAGTGACTGC TGTCTCAGTC** PDL AGTCCCCAGA GGGAAAGACC TTGCCTTAGA AGAAGAGGCG TGTGGGGAAC 951 GGGGGCTCTT GAAGCCAGGT AGCTGGGGAC TATGGTGTCT GCCCTTCCAA BDDB TCACCTCCCT GACCCCTGCT GTCCCATTTT CCCCAGCTGG CCGCATTCCT 1051 CTCTGCTTCT CAGCAGCTGT CCTACTCCCC AGGACGAGTT TTCACTAGAG **JIDI GGCCCACGAT GCCAGGATTC TGATTCATCT TCCTCCCAAG AAAAGCAAAG** 55 1151 CCAAATCAAG ACCACAGATA GGAACCTAAG CACAATGGGG TGCCTGCTTG 12D1 GGCTGGGTCG AAGGCTCTGC TGACTGCTGT CCTTGTCCAT CACCCAATAC 1251 CACCCCAAAC ACAACTCAAC TTCCCACACC ACCATGTCTC TCACCACACC

WO 01/98454 PCT/IB01/02050 BOD TTCTGGGCCT CATTATCTC CACAACTAGA CCGCCATGCC TCACCAACCT
L351 ATGTCCCTGG ACCTCCTGGT GTCTGCCTCT CGGAGTCTGT GCACATCTGC
L4D1 TCACAGTTGA GTGGGGGAAG AAACAGCCAG AATTCAATAC AACAAAGAGC
L451 GGGAGTTAGT ATAGGAATGT CCATCTCATA AGGCTGAGAG CTATTTTTTC
L5D1 CTGTGGCTGC AAATGTCTGA AGCCAGTTAG TTTGATTACC CTGTGCAAAA
L551 CCTTGGACAT ACTTCTGCTA TTAACGCTAT AGGTATTTAT CCGTTTCCAC
L6D1 TGGCTTTTTG TACCCACCGA GCCCCTGAGC CTTGCGTGTG TGTGTGTGGA
L651 AGAGCCTTGT AGAGAACTGC TCCTGTGAGG CAGACAGGAC AGTGAGGTTG
L701 TCACCACTCA GACTTCACCT ATTCAGCATT CTTTCTGATT TCTAGAACTA
L751 TCCACCTCAT TAGGCCTTCT TCCTATCCCC ATCTCTGGCC TCTTGAGCTT
L801 AAGCTTGTAT TGTCCTGGAA TCAGTGGCTT TCTAACCCCC TGCCAGGCTT
L851 TGCCAAAGCA AAAAGACAGA GGCTTTTTTT TTTTTTTTAA AGTTTGGGGT
L901 CTGTCAGGAG ACCAGAGGCT TTTTGAATTC ACTGTGAAGA GAAGAACCCG
L951 AACCTTAAGA CGCCAGATCC CTGAGAGTCT TTCTGGCTGG TTTGAGTCTC
CD01 TCAAATCATG GATTAGGAGT AAAGAAAGAG GCAGGCGCAA TGGCTCATGC
CD51 CTGTAATCCC AGCACTTTGG GAGGCTGAGG TGGGTGGATC ACTTGAGGTC
CD51 CTGTAATCCC AGCACTTTGG GAGGCTGAGG TGGGTGGATC ACTTGAGGTC
CD51 CTGTAATCCC AGCACTTTGG GAGGCTGAGG TGGGTGGATC ACTTGAGGTC
CD51 AGGAGTTTGA GACCAGCCTG GGTAATATGG CAAAACCCCA TCTCTACTAA 5 10 15 2101 AGGAGTTTGA GACCAGCCTG GGTAATATGG CAAAACCCCA TCTCTACTAA . 2151 AAAATACAAA AATTAGCCAG GTATGGTGGT GAACACCTGT AATCCCAGCT 2201 ACTTGGAAGG CTGAGGCATA GGAGTTGCTT GAACCTGGGA GATGGGGGTT 2251 GTAGTGAGCC AAGTTCGTGC CATCGGACTC CAGCCTGGGT GAAGGAGTGA 20 2301 GACCCTGTCT CCAAAAACAA ACAAAAAAGG AGCAGAGAAA GACAGTGGTA 2351 CAGCTAACCT GAACAAGGGA ACTGGGACCG TTGGGCTGAA ACAGTCTTGA 2401 GCCTGGGGTT GACTGGGTTA GAGAAGAACC GGGATGCAAG GAGCTGCCTG 2451 TGACACCTGG CCTGCCCTTT CTCAGCTGCC TCCCCTGCCC TTTCTCAGCT
2501 GCCTCCCCTG CCCTCAGAAG GAAAGGAGAG GGCTCACTTA TCACTTGTGC
2551 CATAGCACCT GGTCTCAAAA TCCTAAAAGC TTTCCTCGCC CTCACTGCCT
2601 TGCTCCACAA GGTCCACTTT CCTGGGTCTT GTGCTGTGCC TTTCCTTGTC
2651 TGCCTCCTGC TGCTTCTGTA ACTGCAGACC CCAGGCCCAA TTGCAAGCCC
2701 TCGCTCAGC TGCTTCTCCA TTGGAATAAA CTCTTGTTTC TCTAAAAAAA 25 30 2751 AAAA **BLAST** Results

35

No BLAST result

40

No Medline entry

45

Peptide information for frame 2

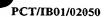
Medline entries

ORF from 41 bp to 883 bp; peptide length: 281

50 Category: putative protein Classification: no clue

Prosite motifs: ASP_PROTEASE (173-184)

J MLFIFPLSLP WRPSCWKESC STGQRQAGRS REDSVTPPPS SPWPTPPAGA
51 MSTKQEARRD EGEARTRGQE AQLRDRAHLS QQRRLKQATQ FLHKDSADLL
101 PLDSLKRLGT SKDLQPRSVI QRRLVEGNPN WLQGEPPRMQ DLIHGQESRR
151 KTSRTEIPAL LVNCKCQDQL LRVAVDTGTQ YNRISAGCLS RLGLEKRVLK



201 ASAGDLAPGP PTQVEQLELQ LGQETVVCSA QVVDAESPEF CLGLQTLLSL

5

BLASTP hits

No BLASTP hits available

10 Alert BLASTP hits for DKFZphamy2_l3gl9, frame 2

PIR:S50646 hypothetical protein YER143w - yeast (Saccharomyces cerevisiae), N = 1, Score = 90, P = 0.26

TREMBL:RNDOLO_1 product: "DNA (cytosine-5-)-methyltransferase";
Rattus
norvegicus mRNA for DNA (cytosine-5-) -methyltransferase; partial
cds.;
N = 1; Score = 81; P = 0.89

20

>PIR:S50646 hypothetical protein YER143w - yeast (Saccharomyces cerevisiae)

Length = 428

25

HSPs:

Score = 90 (13.5 bits), Expect = 3.0e-01, P = 2.6e-01 Identities = 28/112 (25%), Positives = 48/112 (42%)

30

Query: 155 TEIPALLVNCKCQDQLLRVAVDTGTQYNRISAGCLSRLGLEKRVLKASAGD---LAPGPP 211

T++P L +N + + ++ VDTG Q +S + GL + + K G+

+ G

35 Sbict: 199

TQVPMLYINIEINNYPVKAFVDTGAQTTIMSTRLAKKTGLSRMIDKRFIGEARGVGTGKI 258

Query: 212 XXXXXXXXXXXXXXX

CSAQVVDAESPEFCLGLQTLLSLKCCIDLEHGVLRL 263

CS V+D + + +GL L C+DL+

VLR+

Sbjct: 259 IGRIHQAQVKIETQYIPCSFTVLDTDI-

DVLIGLDMLKRHLACVDLKENVLRI 310

45

40

Pedant information for DKFZphamy2_13g19, frame 2

Report for DKFZphamy2_13g19.2

50

ELENGTHI 281

EMM3 31330.97

EpII 8-75

55 EBLOCKSI PRODO49D
EBLOCKSI BPO19216

EPROSITED ASP_PROTEASE

[KW] All_Alpha

	EKM] EKM]	SIGNAL_PEPTIDE 17 LOW_COMPLEXITY 9.96 %
5	SEQ SEG PRD	MLFIFPLSLPWRPSCWKESCSTGQRQAGRSREDSVTPPPSSPWPTPPAGAMSTKQEARRDxxxxxxxxxxx
10	SEQ SEG PRD	EGEARTRG@EA@LRDRAHLS@@RRLK@AT@FLHKDSADLLPLDSLKRLGTSKDL@PRSVI
15	SEQ SEG PRD	QRRLVEGNPNWLQGEPPRMQDLIHGQESRRKTSRTEIPALLVNCKCQDQLLRVAVDTGTQ
20	SEQ SEG PRD	YNRISAGCLSRLGLEKRVLKASAGDLAPGPPTQVEQLELQLGQETVVCSAQVVDAESPEF
20	SEQ SEG PRD	CLGL@TLLSLKCCIDLEHGVLRLKAPFSELPFLPLY@EPG@ cccchhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
25		
		Prosite for DKFZphamy2_13g19.2
30	PS003	L41 173->185 ASP_PROTEASE PDOCOO128
	(No F	Pfam data available for DKFZphamy2_13g19.2)

WO 01/98454

10

DKFZphamy2_14b5

5 group: intracellular transport and trafficing

DKFZphamy2_14b5 encodes a novel 771 amino acid protein which shows 61% identity to the human TYL protein and 48% identity to the human Tic protein.

Both proteins show similarity to Sec7 of Saccharomyces cerevisiae, which takes function in vesicular traficking. The new protein shows also significant similarity to human ARNO3, which is involved in the control of Golgi structure and function.

15 DKFZphamy2_14b5 is predominantly expressed in the cns and germ
cells.

The new protein can find application in diagnosis/therapy of diseases related to vesicular traficking e-g- in synapses of the central nervous system and in studying expression profiles.

similarity to TYL protein (Homo sapiens)

25 Sequenced by EMBL

Locus: /map="445.7 cR from top of Chr5 linkage group"

30 Insert length: 4528 bp
Poly A stretch at pos. 4511, polyadenylation signal at pos. 4489

L CTCGCTCAGC CTCTCCACAT CGCGGCTCCG GCACCTGAAG GGACGCGGGC 51 GGGCGCGGGC AGCTCCGACC GGCGGCGGG GGGCGGGACA GGCAGCCCGG 35 BDB CGGCCTCCG/ TGGCCCCGCC GTGAGAGGCC GGACCCGCGG CGGGGACCAG 151 CAGCGGTCT: ** AGGAGTCCC AGGAGCAGCC AGGACAGGCG GAAGCAGTGG **201 CTGCCATGGALITAGGACAAG CTCTTATCTG CAGTGCCTGA GGAAGGCGAT** 251 GCCACCCGT COCCCGGTCC AGAGCCTGAA GAGGAGCCAG GGGTCCGGAA TGAGGGCC TGAACAGCAG CCTCTGCAGC CCAGGGCACG 40 301 TGGGATGGC 351 AGCGAAGGGO AACCCAGCG GACACTGAGG AACCCACGAA GGACCCAGAT 401 GTGGCCTTCC / GGCCTCAG CCTTGGCCTC TCTCTCACCA ATGGCCTAGC 451 CCTGGGGCCA GACTTGAACA TTCTGGAAGA TTCAGCGGAG TCCAGGCCCT 501 GGAGGGCTGG CGTGCTGGCA GAGGGGGGACA ATGCTTCCAG GAGCCTCTAC 45 551 CCAGATGCTG AGGACCCTCA GCTGGGGTTG GATGGTCCCG GGGAGCCAGA LOD TGTGCGGGAT GGCTTCAGCG CCACGTTTGA GAAGATTCTG GAGTCAGAGC LSD TGCTGCGGGG CACCCAGTAC AGCAGCCTCG ACTCCCTAGA CGGGCTGAGC 701 CTCACGGATG AGAGCGACAG CTGCGTCAGC TTCGAGGCCC CCCTCACACC 751 CCTCATCCAG CAGCGGGCCC GTGACAGCCC TGAGCCAGGG GCTGGGTTGG BD1 GCATTGGGGA CATGGCGTTT GAGGGGGGACA TGGGGGGCAGC TGGTGGTGAT 50 851 GGGGAGCTGG GCAGCCCCCT GCGGCGCTCC ATCTCCAGCA GCCGCTCTGA 901 GAATGTCCTG AGCCGCCTGT CTCTCATGGC CATGCCCAAT GGATTCCATG 951 AAGATGGCCC TCAGGGCCCA GGGGGGGGATG AGGATGATGA TGAGGAGGAC 1001 ACGGACAAGT TGCTGAACTC AGCCAGTGAC CCCAGCCTGA AGGATGGCCT
1051 GTCAGACTCA GACTCTGAGC TCAGCAGCTC GGAGGGGTTG GAGCCTGGTA
1101 GTGCAGACCC TCTGGCCAAC GGGTGCCAGG GGGTCAGTGA AGCTGCTCAT
1151 CGGCTGGCA GCCGTCTCTA CCACCTCGAG GGCTTCCAGC GCTGTGATGT
1201 GGCCCGGCAG CTGGGCAAGA ACAACGAGTT TAGCAGGCTG GTGGCCGGG 55

```
5
       10
     20
    25
    30
  35
  40
                     3251 AGCCTATTTT GGAGCTTCCC CTGTTAGGAA GGATGGCTGC ACCTGGCCCC
3301 CTGGCATTCC TGACGCTCTA GGAGGGAAGG GGGAGGCAGT GCTGGCCTCC
3351 CTTGCCCTGT TTTTCCCTCT TCCAGCTGAC CTGTGACTTA TACTGCTCTT
3401 ACCGATGATA CTTTTGGAAA AAATAGAGCG TGTATGCACC GCCCCGTTTG
3451 TCCCATGGAT ATCCTGGGGT GTGAGTCGGA TGGGACCACG GCCCCGTTTG
3451 TCCCATGGAT ATCCTGGGGT GTGAGTCGGA TGGGACCACG GCCCTGTTTA
3501 TATTTGGGTC TTTATGTTGG TGCTGCCAGG TCTCTGAGCT CCAGAGGTGG
3551 CCTCTTGGAC AGATCTACTG CTATAGGAAT AAAAGACACT CTGTCTCGCA
3601 AATGGCTGCT TGTCAACAAG CCCAAAGATG CTTGTCGGAG GACGGTTATG
3651 GAAGCCCTTA ATTCTTGGTT GTGGGAAAAG GTGGAATGAC AAGTTATTGA
3701 TTGTTTTTCT GTCGCTATTT CTTTCATTTG TCTAGTGAAT CAGAAAGGCT
3751 TAGCCAAGGC CACATCTGGG AAGAGTGGAG AAATTTGCCA CTTGACGATC
3801 ACGGATTAGC TAGCACCTTT AAGCCCTGCA TTTCTCCAAC TGACAAGTGG
3851 GTGGGGGTGA TGGCACCATTC AGTGTGGCTA TGAAGAGCGA ATCCTCTCTA
3901 TTGTTTAAAAT AGATTACTGT AGTTTGGCCA GGAATTTGGC GTCAGTGGTA
3951 ACACACTTAG TTAATAAAAT AAGCCCAGGCT TGCAACTAAG TATCTAACTT
4001 TACAGGCCCA CTCACATTTG AGGCAAGGGG CTATTGAGTA TGTGGAGAGAA
4051 TGTAGTGATT TAAATTCAGA TTATTTAAGT TGGATCAGCT GAAGTGTGTT
4101 TTAGACCCAA ACCATCTGGC CCCTTCGTTT TGCTCAGAGG AAGTAAATGT
45
50
55
```

4151 TCACTTAAAT GAAATTGAAA ACGCCATGTG GCACCACAAA AGAGCTCTCT 4201 GTACTTTCCC CATGCTGCCT CAAAAGTTCT GTGAGTTTCG GGGTCAGTGT 4251 CCCACCCTTC ACTTCCCGAG GGCGGGTGAG TGGAGAGCAG AGCCAGGAGC 4301 TCTGGCAGCT GTGGACAGAT GTGCTTCCTG AGCATGGGTT GTGCCTCCCA 4351 TCAGTAAAAA AATGTTTAGT TCACTTCCTT AATTGTATAA TTATTTTT 4401 GTAAATTATA TACATGTACT ACTGTACTAA AATATTATGT ACATTATAAA 4451 ACATACACAA AAATAGAAAT TTAAAAAAGA TGAGATGAAA ATAAATCTAA 4501 GTCAAAGTTC CAAAAAAAAA AAAAAAAA

10

5

BLAST Results

No BLAST result

15

Medline entries

20 98086482:

Perletti L. Talarico D. Trecca D. Ronchetti D. Fracchiolla NS. Maiolo AT. Neri A.; Identification of a novel gene. PSD. adjacent to

NFKB2/lyt-10, which contains Sec? and pleckstrin-homology

25 domains.

Genomics 46:251-259(1997)

30

Peptide information for frame 2

ORE from 206 bp to 2518 bp; peptide length: 771
35 Category: similarity to known protein
Classification: Cell signaling/communication

1 MEEDKLLSAV PEEGDATRDP GPEPEEPGV RNGMASEGLN SSLCSPGHER
51 RGTPADTEEP TKDPDVAFHG LSLGLSLTNG LALGPDLNIL EDSAESRPWR
40 JOI AGVLAEGDNA SRSLYPDAED PQLGLDGPGE PDVRDGFSAT FEKILESELL
151 RGTQYSSLDS LDGLSLTDES DSCVSFEAPL TPLIQQRARD SPEPGAGLGI
201 GDMAFEGDMG AAGGDGELGS PLRRSISSSR SENVLSRLSL MAMPNGFHED
251 GPQGPGGDED DDEEDTDKLL NSASDPSLKD GLSDSDSELS SSEGLEPGSA
301 DPLANGCQGV SEAAHRLARR LYHLEGFQRC DVARQLGKNN EFSRLVAGEY
45 351 LSFFDFSGLT LDGALRTFLK AFPLMGETQE RERVLTHFSR RYCQCNPDDS
401 TSEDGIHTLT CALMLLNTDL HGHNIGKKMS CQQFIANLDQ LNDGQDFAKD
451 LLKTLYNSIK NEKLEWAIDE DELRKSLSEL VDDKFGTGTK KVTRILDGGN
501 PFLDVPQALS ATTYKHGVLT RKTHADMDGK RTPRGRRGWK KFYAVLKGTI
551 LYLQKDEYRP DKALSEGDLK NAIRVHHALA TRASDYSKKS NVLKLKTADW
50 LOI RVFLFQAPSK EEMLSWILRI NLVAAIFSAP AFPAAVSSMK KFCRPLLPSC
51 TTRLCQEEQL RSHENKLRQL TAELAEHRCH PVERGIKSKE AEEYRLKEHY
701 LTFEKSRYET YIHLLAMKIK VGSDDLERIE ARLATLEGDD PSLRKTHSSP
751 ALSQGHVTGS KTTKDATGPD T

55

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_14b5, frame 2

5 PIR: 601205 TYL protein - human, N = 2, Score = 1421, P = 8.6e-150

TREMBL:AB023159_1 gene: "KIAA0942"; product: "KIAA0942 protein"; Homo

sapiens mRNA for KIAAD942 protein, partial cds., N = 1, Score =

10 1251 P = 2.3e-127

TREMBL:Ub3127_1 gene: "TIC"; product: "Tic"; Human SEC7 homolog

15 (TIC) mRNA, complete cds., N = 1, Score = 1050, P = 4.6e-106

>PIR:GD1205 TYL protein - human Length = 645

20

HSPs:

Score = 1421 (213-2 bits), Expect = 8.6e-150, Sum P(2) = 8.6e-150

25 Identities = 280/452 (61%), Positives = 336/452 (74%)

Query: 301
DPLANGCQGVSEAAHRLARRLYHLEGFQRCDVARQLGKNNEFSRLVAGEYLSFFDFSGLT 360
D L+NG + EAA RLA+RLY L+GF++ DVAR LGKNN+FS+LVAGEYL

30 FF F+G+T
Sbjct: 166
DTLSNGQKADLEAAQRLAKRLYRLDGFRKADVARHLGKNNDFSKLVAGEYLKFFVFTGMT 225

Query: 361

35 LDGALRTFLKAFPLMGETQERERVLTHFSRRYCQCNPDDSTSEDGIHTLTCALMLLNTDL 420
LD ALR FLK LMGETQERERVL HFS+RY QCNP+ +SEDG
HTLTCALMLLNTDL
Sbjct: 22b
LDQALRVFLKELALMGETQERERVLAHFSQRYFQCNPEALSSEDGAHTLTCALMLLNTDL 285

40

Query: 421
HGHNIGKKMSCQQFIANLDQLNDGQDFAKDLLKTLYNSIKNEKLEWAIDEDELRKSLSEL 480
HGHNIGK+M+C FI NL+ LNDG DF ++LLK

LY+SIKNEKL+WAIDE+ELR+SLSEL

45 Sbjct: 286
HGHNIGKRMTCGDFIGNLEGLNDGGDFPRELLKALYSSIKNEKLQWAIDEEELRRSLSEL 345

Query: 481 VDDKFGTGTKKVTRIL---DGGNPFLDVPQALSATTYKHGVLTRKTHADMDGKRTPRGR 536

50 D K + RI G +PFLD+ A YKHG L RK HAD D ++TPRG+
Sbjct: 346 ADPN---PKVIKRISGSGSGSSPFLDLTPEPGAAVYKHGALVRKVHADPDCRKTPRGK 401

55 Query: 537
RGWKKFYAVLKGTILYLQKDEYRPDKALSEGDLKNAIRVHHALATRASDYSKKSNVLKLK 596
RGWK F+ +LKG ILYLQK+EY+P KALSE +LKNAI
+HHALATRASDYSK+ +V L+

Sbjct: 402

RGWKSFHGILKGMILYLQKEEYKPGKALSETELKNAISIHHALATRASDYSKRPHVFYLR 461

Query: 597

RPLLPS TRL @ Sbict: 462

TADWRVFLFQAPSLEQMQSWITRINVVAAMFSAPPFPAAVSSQKKFSRPLLPSAATRLSQ 521

10

Query: 657
EEQLRSHENKLRQLTAELAEHRCHPVERGIKSKEAEEYRLKEHYLTFEKSRYETYIHLLA 716
EEQ+R+HE KL+ + +EL EHR + + + KEAEE R KE YL FEKSRY

15 Sbjct: 522

EEQVRTHEAKLKAMASELREHRAAQLGKKGRGKEAEEQRQKEAYLEFEKSRYSTYAALLR 581

Query: 717 MKIKVGSDDLERIEARLATLEGDDPSLRKTHSSPAL 752 +K+K GS++L+ +EA LA + L +HSSP+L

20 Sbjct: 582 VKLKAGSEELDAVEAALAQAGSTEDGLPPSHSSPSL 617

Score = 63 (9.5 bits), Expect = 8.6e-150, Sum P(2) = 8.6e-150

Identities = 19/64 (29%), Positives = 23/64 (35%)

D D FS FE ILES +GT Y

+FE P P

Sbjct: 18

30 DGPDSFSCVFEAILESHRAKGTSYTSLASLEALASPGPTQSPFFTFELPPQPPAPRPDPP 77

Query: 191 SPEP 194

+P P

Sbjct: 78 APAP 81

35

Pedant information for DKFZphamy2_14b5, frame 2

40

Report for DKFZphamy2_14b5.2

ELENGTHD 771

EMW3 84660.55

45 EpII 5.04

CHOMOLD PIR:GOL205 TYL protein - human le-158

EFUNCATI 30.09 organization of intracellular transport vesicles
ES. cerevisiae, YDR170cl 5e-22

[FUNCAT] 30.08 organization of golgi [S. cerevisiae, YDR]70c]

50 5e-22

EFUNCATI 30.03 organization of cytoplasm **ES**· cerevisiae.

YDR170c1 5e-22

EFUNCATD D8-D7 vesicular transport (golgi network, etc.)
ES.

cerevisiae, YDR170c1 5e-22

55 [FUNCAT] 99 unclassified proteins [S. cerevisiae, YPR095c]

4e-04

EBFOCKZI BF07534B

EBLOCKSI BPD2373F

EBLOCKSI PROOL55C EBLOCKSD PROJUBBE EBFOCKZ] bb005548 EBFOCKZI Bb05P4PD EBLOCKSI PRODERLA EBFOCKZ3 DWO7324W **EBLOCKZI** - PFD13h9B RBLOCKSI PFO1369A ESCOPI . dlbtn_ 2.41.1.1.2 beta-spectrin Emouse (Mus musculus) brain le-39 10 [PIRKW] transmembrane protein le-20 ESUPFAMD Caenorhabditis elegans KObH7.4 protein 7e-24
ESUPFAMD pleckstrin repeat homology 7e-24 EPFAMI PH (pleckstrin homology) domain 15 EKWI . Irregular EKWI ΒD LOW_COMPLEXITY 18-42 % **EKWI** 20 SEQ MEEDKLLSAVPEEGDATRDPGPEPEEEPGVRNGMASEGLNSSLCSPGHERRGTPADTEEP SEG ····· lbtn-25 SEQ TKDPDVAFHGLZLGLZLTNGLALGPDLNILEDSAESRPWRAGVLAEGDNASRSLYPDAED SEG lbtn-30 SEQ PQLGLDGPGEPDVRDGFSATFEKILESELLRGTQYSSLDSLDGLSLTDESDSCVSFEAPL SEG lbtn-35 SEQ-- TPLIQQRARDSPEPGAGLGIGDMAFEGDMGAAGGDGELGSPLRRSISSSRSENVLSRLSL SEG lbtn-SEQ MAMPNGFHEDGPQGPGGDEDDDEEDTDKLLNSASDPSLKDGLSDSDSELSSSEGLEPGSA 40 SEG lbtn-45 SEQ DPLANGCQGVSEAAHRLARRLYHLEGFQRCDVARQLGKNNEFSRLVAGEYLSFFDFSGLT SEG lbtn-SEQ LDGALRTFLKAFPLMGETQERERVLTHFSRRYCQCNPDDSTSEDGIHTLTCALMLLNTDL 50 SEG lbtn-55 HGHNIGKKMSCQQFIANLDQLNDGQDFAKDLLKTLYNSIKNEKLEWAIDEDELRKSLSEL SEG lbtn-

PCT/IB01/02050

WO 01/98454

	SEQ VDDKFGTGTKKVTRILDGGNPFLDVPQALSATTYKHGVLTRKTHADMDGKRTPRGRRGWK
5	SEGEEEEEEEETTTEET TTTCEE
	SEQ KFYAVLKGTILYLQKDEYRPDKALSEGDLKNAIRVHHALATRASDYSKKSNVLKLKTADW SEG
10	lbtn- EEEEEETTEEEEECCHHHHHHHCCBTTT- TCCEETTTTEEEEETTTTTTTTTTTTEEEEETTTT
	SEQ RVFLFQAPSKEEMLSWILRINLVAAIFSAPAFPAAVSSMKKFCRPLLPSCTTRLCQEEQL SEGxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
15	СЕЕЕЕСССНИННИНИННИН
	SEQ RSHENKLRQLTAELAEHRCHPVERGIKSKEAEEYRLKEHYLTFEKSRYETYIHLLAMKIK SEG
20	•
0.5	SEQ VGSDDLERIEARLATLEGDDPSLRKTHSSPALSQGHVTGSKTTKDATGPDT SEG
25	
	(No Prosite data available for DKFZphamy2_14b5.2)
30	Pfam for DKFZphamy2_14b5-2
30	Pfam for DKFZphamy2_14b5.2 HMM_NAME PH (pleckstrin homology) domain
30 35	HMM_NAME PH (pleckstrin homology) domain HMM *dvIREGUMyKUgswrkstg······nUqrRUFvLrndpnrLiYYkddk + ++G + +++ + ++
	HMM_NAME PH (pleckstrin homology) domain HMM *dvIREGWMyKWgswrkstg·····nWqrRWFvLrndpnrLiYYkddk
35	HMM_NAME PH (pleckstrin homology) domain HMM *dvIREGWMyKWgswrkstgnWqrRWFvLrndpnrLiYYkddk + ++G + +++ + ++
35	HMM_NAME PH (pleckstrin homology) domain HMM *dvIREGWMyKWgswrkstgnWqrRWFvLrndpnrLiYYkddk + ++G + +++ + ++
35 40 45	HMM_NAME PH (pleckstrin homology) domain HMM *dvIREGWMyKWgswrkstgnWqrRWFvLrndpnrLiYYkddk + ++G + +++ + ++
35 40	HMM_NAME PH (pleckstrin homology) domain HMM *dvIREGWMyKWgswrkstgnWqrRWFvLrndpnrLiYYkddk + ++G + +++ + ++

10

15

20

30

DKFZphamy2_14mlb

5 group: transcription factors

DKFZphamy2_14m16.pl encodes a novel 252 amino acid protein with similarity to the homeotic protein emx2 of man, mouse and zebra fish as well as to the gene "empty spiracles" of Drosophila melanogaster.

Homoeobox genes are known to play important roles in developmental processes. In zebrafish emx2 mRNAs are found in the dorsal telencephalon, parts of the diencephalon and the otocyst. The human homologue Emx2 appears to be already expressed in 8.5 day embryos. It is also expressed in the presumptive cerebral cortex, olfactory bulbs, in some neuroectodermal areas in embryonic head including olfactory placodes in earlier stages and olfactory epithelia later in development. Mutants of the D. melanogaster gene "mempty spiracles" display spiracles devoid of filzkorper, no antenna and an open head.

The new protein can find application in modulating the expression of genes controlled by this transcription factor and modulation of neuronal development.

strong similarity to homeotic protein emx2 (Homo sapiens)

perhaps differential splicing

Sequenced by EMBL

Locus: /chromosome="l0"

35 Insert length: 2416 bp Poly A stretch at pos. 2378 polyadenylation signal at pos. 2373

I GAAAAAAA GAAAAAAA GAAAAAAAT TACCCCAATC CACGCCTGCA 51 AATTCTTCTG GAAGGATTTT CCCCCCTCTC TTCAGGTTGG GCGCGTTTGG
101 TGCAAGATTC TCGGGATCCT CGGCTTTGCC TCTCCCTCTC CCTCCCCCT
151 CCTTTCCTTT TTCCTTTCCT TTCCTTTCTT TCTTCCTTTC CTTCCCCCA
201 CCCCACCC CACCCCAAAC AAACGAGTCC CCAATTCTCG TCCGTCCTCG 40 251 CCGCGGGCAG CGGGCGGCGG AGGCAGCGTG CGGCGGTCGC CAGGAGCTGG
301 GAGCCCAGGG CGCCCGCTCC TCGGCGCAGC ATGTTCCAGC CGGCGCCCAA
351 GCGCTGCTTC ACCATCGAGT CGCTGGTGGC CAAGGACAGT CCCCTGCCCG
401 CCTCGCGCTC CGAGGACCCC ATCCGTCCCG CGGCACTCAG CTACGCTAAC 45 451 TCCAGCCCCA TAAATCCGTT CCTCAACGGC TTCCACTCGG CCGCCGCCGC 501 CGCCGCCGGT AGGGGCGTCT ACTCCAACCC GGACTTGGTG TTCGCCGAGG 551 CGGTCTCGCA CCCGCCCAAC CCCGCCGTGC CAGTGCACCC GGTGCCGCCG 50 LOT CCGCACGCCC TGGCCGCCCA CCCCCTACCC TCCTCGCACT CGCCACACCC 651 CCTATTCGCC TCGCAGCAGC GGGATCCGTC CACCTTCTAC CCCTGGCTCA 701 TCCACCGCTA CCGATATCTG GGTCATCGCT TCCAAGGGAA CGACACTAGC 751 CCCGAGAGTT TCCTTTTGCA CAACGCGCTG GCCCGAAAGC CCAAGCGGAT 55 BOL CCGAACCGCC TTCTCCCCGT CCCAGCTTCT AAGGCTGGAA CACGCCTTTG ASI AGAAGAATCA CTACGTGGTG GGCGCCGAAA GGAAGCAGCT GGCACACAGC PDL CTCAGCCTCA CGGAAACTCA GGTAAAAGTA TGGTTTCAGA ACCGAAGAAC 951 AAAGTTCAAA AGGCAGAAGC TGGAGGAAGA AGGCTCAGAT TCGCAACAAA

LODL AGAAAAAGG GACGCACCAT ATTAACCGGT GGAGAATCGC CACCAAGCAG 1051 GCGAGTCCGG AGGAAATAGA CGTGACCTCA GATGATTAAA AACATAAACC BIDI TAACCCGACA GAAACGGACA ACATGGAGCA AAAGAGACAG GGAGAGGTGG 1151 AGAAGGAAAA AACCCTACAA AACAAAAACA AACCGCATAC ACGTTCACCG 1201 AGAAAGGAAG AGGGAATCGG AGGGAGCAGC GGAATGCGGC GAAGACTCTG
1251 GACAGCGAGG GCACAGGGTC CCAAACCGAG GCCGCGCCAA GATGGCAGAG 5 1301 GATGGAGGCT CCTTCATCAA CAAGCGACCC TCGTCTAAAG AGGCAGCTGA 1351 GTGAGAGACA CAGAGAGAG GAGAAAGAG GAGGGAGAGA GAGAAAGAGA 1401 GAGAAAGAGA GAGAGAGAGA GAGAGAAAGC TGAACGTGCA CTCTGACAAG 1451 GGGAGCTGTC AATCAAACAC CAAACCGGGG AGACAAGATG ATTGGCAGGT 1501 ATTCCGTTTA TCACAGTCCA CTTAAAAAAT GATGATGATG ATAAAAACCA 1551 CGACCCAACC AGGCACAGGA CTTTTTTGTT TTTTGCACTT CGCTGTGTTT 3603 CCCCCCATC TTTAAAAATA ATTAGTAATA AAAAACAAAA ATTCCATATC 1651 TAGCCCCATC CCACACCTGT TTCAAATCCT TGAAATGCAT GTAGCAGTTG 1701 TTGGGCGAAT GGTGTTTAAA GACCGAAAAT GAATTGTAAT TTTCTTTTCC 15 1751 TTTTAAAGAC AGGTTCTGTG TGCTTTTTAT TTTGATTTTT TTTCCCAAGA 1801 AATGTGCAGT CTGTAAACAC TTTTTGATAC CTTCTGATGT CAAAGTGATT 1851 GTGCAAGCTA AATGAAGTAG GCTCAGCGAT AGTGGTCCTC TTACAGAGAA 1901 ACGGGGAGCA GGACGACGGG GGGGCTGGGG GTGGCGGGG AGGGTGCCCA 1951 CAAAAAGAAT CAGGACTTGT ACTGGGAAAA AAACCCCTAA ATTAATTATA 20 2001 TTTCTTGGAC ATTCCCTTTC CTAACATCCT GAGGCTTAAA ACCCTGATGC 2051 AAACTTCTCC TTTCAGTGGT TGGAGAAATT GGCCGAGTTC AACCATTCAC 2101 TGCAATGCCT ATTCCAAACT TTAAATCTAT CTATTGCAAA ACCTGAAGGA 2151 CTGTAGTTAG CGGGGATGAT GTTAAGTGTG GCCAAGCGCA CGGCGGCAAG 25 2201 TTTTCAAGCA CTGAGTTTCT ATTCCAAGAT CATAGACTTA CTAAAGAGAG 2251 TGACAAATGC TTCCTTAATG TCTTCTATAC CAGAATGTAA ATATTTTTGT 2301 GTTTTGTGTT AATTTGTTAG AATTCTAACA CACTATATAC TTCCAAGAAG 2351 TATGTCAATG TCAATATTTT GTCAATAAAG ATTTATCAAT ATGCCCTCAC AAAAAA AAAAAA COPS 30

BLAST alert EMBL/EMBLNEW

35 EMBLNEW:ALl33353 Human DNA sequence *** SEQUENCING IN PROGRESS *** from clone RPl1-483Fl1: N = 2, Score = 3108, P = 5.3e-134

EMBL:HSEMX2 H-sapiens EMX2 mRNA; N = 1, Score = 2385, P = 5.1e-40 101

Medline entries

45 92331606:
Simeone A. Gulisano M. Acampora D. Stornaiuolo A. Rambaldi M. Boncinelli E.;
Two vertebrate homeobox genes related to the Drosophila empty spiracles gene are expressed in the embryonic cerebral cortex.
50 EMBO J
1992 Jul;11(7):2541-50

55

ORF from 331 bp to 1086 bp; peptide length: 252

Category: questionable ORF

Classification: unset

Prosite motifs: HOMEOBOX_1 (187-210)

5

10

1 MFQPAPKRCF TIESLVAKDS PLPASRSEDP IRPAALSYAN SSPINPFLNG 51 FHSAAAAAAG RGVYSNPDLV FAEAVSHPPN PAVPVHPVPP PHALAAHPLP 101 SSHSPHPLFA SQQRDPSTFY PWLIHRYRYL GHRFQGNDTS PESFLLHNAL 151 ARKPKRIRTA FSPSQLLRLE HAFEKNHYVV GAERKQLAHS LSLTETQVKV 201 WFQNRRTKFK RQKLEEEGSD SQQKKKGTHH INRWRIATKQ ASPEEIDVTS 251 DD

15 Alert BLASTP hits for DKFZphamy2_14mlb, frame 1

PIR:I51737 homeotic protein emx2 - zebra fish; N = 2, Score = 753 P = le-105

20

PIR:S22722 homeotic protein emx2 - human (fragment); N = la Score $763_{1} P = 1.3e-75$

TREMBL: OLA132403_1 gene: "emx2"; product: "Emx2 protein"; 25 latipes mRNA for Emx2 protein, partial; N = 2, Score = 513, P = 4 · 5e-72

30

>PIR:S22722 homeotic protein emx2 - human (fragment) Length = 158

HSPs:

35

45

Score = 763 (114.5 bits), Expect = 1.3e-75, P = 1.3e-75Identities = 144/144 (100%), Positives = 144/144 (100%)

Query: 109

FASQQRDPSTFYPWLIHRYRYLGHRFQGNDTSPESFLLHNALARKPKRIRTAFSPSQLLR 168 40

FASQQRDPSTFYPWLIHRYRYLGHRFQGNDTSPESFLLHNALARKPKRIRTAFSPSQLLR Sbjct: . 1.5 FASQQRDPSTFYPWLIHRYRYLGHRFQGNDTSPESFLLHNALARKPKRIRTAFSPSQLLR 74

Query: **169** LEHAFEKNHYVVGAERKQLAHSLSLTETQVKVWFQNRRTKFKRQKLEEEGSDSQQKKKGT 228

LEHAFEKNHYVVGAERKQLAHSLSLTETQVKVWFQNRRTKFKRQKLEEEGSDSQQKKKGT 50 Sbict: 75 LEHAFEKNHYVVGAERKQLAHSLSLTETQVKVWFQNRRTKFKRQKLEEEGSDSQQKKKGT 134

Query: 229 HHINRWRIATKQASPEEIDVTSDD 252 HHINRURIATKQASPEEIDVTSDD

55 Sbjct: 135 HHINRWRIATKQASPEEIDVTSDD 158

Pedant information for DKFZphamy2_14mlb, frame 1

Report for DKFZphamy2_14m16.1

```
5
    ELENGTH
                     362
                     40749.28
    [[q]
                     10.51
                     PIR:I51737 homeotic protein emx2 - zebra fish le-
    EHOMOL
    113
10
    EFUNCATI
                     30.10 nuclear organization
                                                     ES - cerevisiae -
    YML027w3 5e-05
    EFUNCATI
                    04-99 other transcription activities
                                                              EZ-
    cerevisiae, YMLO27w1 5e-05
    EFUNCATI
                     03.07 pheromone response, mating-type
15 determination, sex-specific proteins
                                               EZ.
    cerevisiae, YCRO97w3 5e-84
    EFUNCATE
                     04.05.01.04 transcriptional control
                                                              EZ-
    cerevisiae, YDL106c1 7e-04
                     01.04.04 regulation of phosphate utilization
    EFUNCATE
20
    ES. cerevisiae, YDLLObcl 7e-04
    EFUNCATE
                    01-03-13 regulation of nucleotide metabolism
    ES. cerevisiae, YDLlObcl 7e-04
    EBF0CK21
                    PR00049D
    EBF0CK23
                    PRODUDAH
25
    EBFOCKZ
                    PRDO487F
    EBFOCKZ
                    PR007966
    EBF.OCK21
                    BL00035C
    EBFOCKZ
                    BL00027 'Homeobox' domain proteins
    EBFOCK2
                    PR00026A
30
                    BL00032C
    EBFOCK21
                    BLDDD32B 'Homeobox' antennapedia-type protein
    EBFOCK2
                    dlau7bl 1.4.1.1.6 Pit-1 POU homeodomain Pit-1
    EZCOPI
    Pit-1
           CRat (Rattu 5e-16
    EZCOPI
                    dlyrna_ 1.4.1.1.2 mating type protein Al
35
    Homeodomain mat alpha 2e-15
    ESCOPI
                    dlenh__ 1.4.1.1.1 engrailed Homeodomain
    [(Drosophila melanogaster 2e-13
    CPIRKWI
                    nucleus le-67
    EPIRKWI
                    heart 3e-10
40
    [PIRKW]
                    DNA binding le-67
    EPIRKWJ
                    leukemia 3e-15
    CPIRKUJ
                    alternative splicing le-10
    EPIRKWI
                    proto-oncogene 3e-15
    EPIRKWI
                    transcription factor be-11
45
    EPIRKWI
                    embryo 9e-12
    EPIRKU
                    transcription regulation Le-67
    EPIRKUJ
                    homeobox le-67.
    ESUPFAMI
                    homeobox homology le-67
    ESUPFAMI
                    homeotic protein Hox A5 7e-10
50
    ESUPFAMI
                    homeotic protein Hox B3 3e-10
    ESUPFAMJ
                    homeotic protein Hox B2 3e-11
    ESUPFAMD
                    homeotic protein Hox Bl 7e-ll
    ESUPFAMD
                    unassigned homeobox proteins le-67
    ESUPFAMD
                    homeotic protein goosecoid 4e-10
55
    ESUPFAMI
                    homeotic protein Hox D4 9e-12
    EPROSITE
                    HOMEOBOX_1
    EPFAMI
                    Homeobox domain
    EKW1
                    Irregular
```

	WO 01/98454		PCT/IB01/02050
	EKMJ EKMJ	3D LOW_COMPLEXITY 25.6	.9 %
5	SEQ		
		PRLQILLEGFSPLSSGWARLVQDSRD	PRLCLSLSLPPPFLFPFLSFL
10	lfilA		
10	SEQ	••••••••••••••••••••••	•••••
	SZEPSPHPHPHPK SEG	QTSPQFSSVLAAGSGRRRQRAAVARS	
15	xxxxxxxxxxxx. lfjlA		(
	ZEQ	•••••••••••••••••••••••••••••••••••••••	•••••••
20	SEG	ASRSEDPIRPAALSYANSSPINPFLN	
	lf ilA	• • • • • • • • • • • • • • • • • • • •	•
25	SEQ	PJAHASHSSAJAHAATAHAAAAAAAA	
30	····xxxxxxxxxxx	××××××××××××××××××××××××××××××××××××××	
35	SEQ GHRFQGNDTSPESI SEG	FLLHNALARKPKRIRTAFSPSQLLRL	.EHAFEKNHYVVGAERKQLAHS
33		• • • • • • • • • • • • • • • • • • • •	•••••
		СССССССССНННННН	ннинниттттснинниннин
40	ZEG CZLTETQVKVWFQI ZEQ	NRRTKFKRØKLEEEGSDSØØKKKGTH	
45	Ъf ј1А НСССННННННННН		

Pfam for DKFZphamy2_14mlb-1

Prosite for DKFZphamy2_14ml6.1

PD0C00027

HOMEOBOX_1

SEQ SEG lfjla

PS00027

297->321

50

55

HMM_NAME

Homeobox domain

MMH

5 *RRRPRTtFTre@LdELEREFHfNrYPTRqRREELA@mLNLTER@VKIWF

+R RT+F+ +QL++LE +F+ N+Y+ ++R

+LA++L+LTE+QVK+UF

Query 264

PKRIRTAFSPSQLLRLEHAFEKNHYVVGAERKQLAHSLSLTETQVKVWF 375

10 HMM

@NRRMKWKRMH*

QNRR+K KR+

*Q*uery

313 QNRRTKFKRQK

323

15

WO 01/98454 DKFZphamy2_16e14

5 group: amygdala derived

DKFZphamy2_16e14-p3 encodes a novel 328 amino acid proteins similar to carbonic anhydrase-related proteins.

- A similar cDNA encoding a protein of the same length was identified in sheep. This protein shows a strong signal sequence, which indicates that it is a secreted protein. The new protein belongs to a protein family, which was designated carbonic anhydrase-related protein XI (CA-RP XI), encoded by CALL (human)
- and Carll (mouse, rat). Despite potentially inactivating changes in the active-site residues, CA-RP XI is evolving very slowly in mammals, a property indicative of an important function, which has also been observed in the two other "acatalytic" CA isoforms, CA-RP VIII and CA-RP X.
- 20 No informative BLAST results: No predictive prosite, pfam or SCOP motife.

The new protein can find application in studying the expression profile of amygdala-specific genes.

similarity to carbonic anhydrase-related protein (Homo sapiens)

ESTs ending at appr. 1800 have polyA-signal

30 Sequenced by EMBL

25

Locus: /map="17q24; 5.13cR from GATA41CD5"

35 Insert length: 2267 bp
Poly A stretch at pos. 2252, polyadenylation signal at pos. 2231

L GGATGGAAAT AGTCTGGGAG GTGCTTTTTC TTCTTCAAGC CAATTTCATC 40 51 GTCTGCATAT CAGCTCAACA GAATTCACCA AAAATCCATG AAGGCTGGTG **JDL GGCATACAAG GAGGTGGTCC AGGGAAGCTT TGTTCCAGTT CCTTCTTTCT** 151 GGGGATTGGT GAACTCAGCT TGGAATCTTT GCTCTGTGGG GAAACGGCAG 201 TCGCCAGTCA ACATAGAGAC CAGTCACATG ATCTTCGACC CCTTTCTGAC 251 ACCTCTTCGC ATCAACACGG GGGGCAGGAA GGTCAGTGGG ACCATGTACA BOL ACACTGGAAG ACACGTATCC CTTCGCCTGG ACAGGAGCA CTTGGTCAAC 45 351 ATATCTGGAG GGCCCATGAC ATACAGCCAC CGGCTGGAGG AGATCCGACT 401 ACACTTTGGG AGTGAGGACA GCCAAGGGTC GGAGCACCTC CTCAATGGAC 451 AGGCCTTCTC TGGGGAGGTG CAGCTCATCC ACTATAACCA TGAGCTATAT 501 ACGAATGTCA CAGAAGCTGC AAAGAGTCCA AATGGATTGG TGGTAGTTTC 50 551 TATATTTATA AAAGTTTCTG ATTCATCAAA CCCATTTCTT AATCGAATGC LOT TCAACAGAGA TACTATCACA AGAATAACAT ATAAAAATGA TGCATATTTA 651 CTACAGGGGC TTAATATAGA GGAACTATAT CCAGAGACCT CTAGTTTCAT 701 CACTTATGAT GGGTCGATGA CTATCCCACC CTGCTATGAG ACAGCAAGTT 751 GGATCATAAT GAACAAACCT GTCTATATAA CCAGGATGCA GATGCATTCC BOL TTGCGCCTGC TCAGCCAGAA CCAGCCATCT CAGATCTTTC TGAGCATGAG
BSL TGACAACTTC AGGCCTGTCC AGCCACTCAA CAACCGCTGC ATCCGCACCA
POL ATATCAACTT CAGTTTACAG GGGAAGGACT GTCCAAACAA CCGAGCCCAG
PSL AAGCTTCAGT ATAGAGTAAA TGAATGGCTC CTCAAGTAGG GAACAAAGCC 55

	WO 01/98454					PCT/IB01/02050
	7007		CCACCTCAGT		AACTGTGAAT.	TGACGTAACC
	1051		CCCTTCTTGC	TTCTCTCTCC	TTCTTTCCCC	CAAGCCTCAT
	1101	TCATTCTTGG	GATTGGCCCT	TTCTTCATGA	AAAGTGTCTG	CAAAACCATG
	1151	GCAGAGGAAT	ACATCTCTCA	CACATACTCA	CAAACACACA	CACAAGCACT
5	1501	TGCACATACA	TACAAACACA	TGCAAACATA	CCTACACACA	CACACACTCT
	7527	TACAACCTCC	ATCATGGGAA	GTCAAGTTTC	AGAAACAAAA	GTCTCATTCA
	7307	TAAGAGGTCT	TAGAAGAAAA	TAACCAGTTA	ACCTGATTTC	AATTTTGATA
	1351	CCGTTTTCCT	GAACTAATAA	ATCTACCCAA	TGAGACTTTT	CAGCCTTTGT
	1401	ACATACAAAA	TTCTTCCAAA	AGAGAGAGGA	GAAAATACAG	CTCTGATGGC
10 .	1451	ATCAAACGGA	CTTTGCATCA	AGTAATTTCA	GATAGTGTCC	TAGGATCCTT
	1501	TGAGGGTGCT	GGTAGCAGGT	GAGCAGGACA	AAGTTGACCA	AGGACACTTA
		TTTCTAGATT	ATGATTCTTC	TGTTTACTCA	ACAATTTACA	AAGAAAAAA
		GGACAGACAT	TGAAGAGCTA	CACATTGTAT	ATATATCACC	ACAGACTATA.
	3 P 2 J	AGGAAATGGA	ATTATTTCCC	TCTTTGTCAC	ATATCTGTAG	TAGGATTTGC
15	1701	CAAGATCAGA	AATGATCCAT	TTGCTGTTTC	TTGTTTTCCA	AAGGTCATAC
	1751	ATTGTGTTTG	GTTATTGTTA	CCAGCTCAAT	AAATGTGTTT	AACGAGTTAA
	7907	TTTCATTTTT	CTGGCTTTGG	TCTGTTCTCC	TTCCTTACAG	GCTAAGCCCT
	1851	GGCTCCATGC	AACTGCATTC	TTTGATTTCA	CTTGTTCCTT	CATCTACATG
	1901	TTTTGTTCAT	TTGCAGCCAG	TTTTTACTGA	GTTTGTGGCA	ATCAGGAATG
20	1951	CATTTGCTAA	GCAAGTATGA	CTTTAATTCC	ACTCCATGGC	TCAATCATTC
	5007	ACATGAGGTG	AGCTTCAGCC	TGAGATAGCA	GGCGACAGAC	TTCTTGCGTT
	5027	TCAAAACTGC	CATGCCCCC	TGTGATGCTC	CCGTGAAGGA	ATGCACTTTG
	5707	CCTTGTAAGT	TCCTGGGAAA	GGGGTATGTT	TTCTCTCCAG	GTGCAGCCAG
~~	2727	ATCTCACAAA	GTACAAAACG	AATGCCTTTC	TTTTCTTGTT	TATAATGGTC
25	5507	ACTCACTGTG	TTTGGTTACT	GTCAAGAAAT	CAATAAATGT	GTTTAACAAG
	2251	TCAAAAAAAA	AAAAAA			

BLAST alert EMBL/EMBLNEW

30

EMBL:AFD64854 Homo sapiens map 17q24; 5-13cR from GATA41CD5 repeat region, complete sequence.; N = 2, Score = 8784, P = 0

3*5*

EMBLNEW:ACOO5883 Homo sapiens chromosome 17 clone RP11-958E11 map 17.
WORKING DRAFT SEQUENCE, 2 ordered pieces.; N = 3, Score = 6260, P = 0

40

Medline entries

9097349:

- 45 Lovejoy DA: Hewett-Emmett D: Porter (A: Cepoi D: Sheffield A: Vale WW: Tashian RE: Evolutionarily conserved: "acatalytic" carbonic anhydrase-related protein XI contains a sequence motif present in the neuropeptide sauvagine:
- 50 the human CA-RP XI gene (CAll) is embedded between the secretor gene cluster and the
- 55 DBP gene at 19q13.3. Genomics 1998 Dec 15:54(3):484-9

Peptide information for frame 3

5 ORF from 0 bp to 986 bp; peptide length: 329 Category: similarity to known protein Classification: unclassified

1 MEIVWEVLFL LQANFIVCIS AQQNSPKIHE GWWAYKEVVQ GSFVPVPSFW
10 51 GLVNSAWNLC SVGKRQSPVN IETSHMIFDP FLTPLRINTG GRKVSGTMYN
101 TGRHVSLRLD KEHLVNISGG PMTYSHRLEE IRLHFGSEDS QGSEHLLNGQ
151 AFSGEVQLIH YNHELYTNVT EAAKSPNGLV VVSIFIKVSD SZNPFLNRML
201 NRDTITRITY KNDAYLLQGL NIEELYPETS SFITYDGSMT IPPCYETASW
251 IIMNKPVYIT RMQMHSLRLL SQNQPSQIFL SMSDNFRPVQ PLNNRCIRTN
15 301 INFSLQGKDC PNNRAQKLQY RVNEWLLK

Alert BLASTP hits for DKFZphamy2_16e14, frame 3

20 PIR:JED375 carbonic anhydrase-related protein - human; N = la Score = 937; P = 4.6e-94

SWISSNEW:CAHB_SHEEP CARBONIC ANHYDRASE-RELATED PROTEIN 2

PRECURSOR
(CARP 2) (CA-RP II) (CA-XI).; N = 1, Score = 935, P = 7.5e-94

>PIR:JEO375 carbonic anhydrase-related protein - human Length = 328

HSPs:

30

50

Score = 937 (140.6 bits), Expect = 4.6e-94, P = 4.6e-94 35 Identities = 169/287 (58%), Positives = 223/287 (77%)

Query: 30
EGWWAYKEVVQGSFVPVPSFWGLVNSAWNLCSVGKRQSPVNIETSHMIFDPFLTPLRINT 89
E WW+YK+ +QG+FVP P FWGLVN+AW+LC+VGKRQSPV++E

40 +++PFL PLR++T
Sbjct: 32
EDWWZYKDNLQGNFVPGPPFWGLVNAAWSLCAVGKRQSPVDVEVKRVLYDPFLPPLRLST 71

Query:9045GGRKVSGTMYNTGRHVSLRLDKEHLVNISGGPMTYSHRLEEIRLHFGSEDSQGSEHLLNG 149GG K+ GT+YNTGRHVS+VN+SGGP+ YSHRL E+RL FG+ DGSEH +NGSEH +NSbjct:92

GGEKLRGTLYNTGRHVSFLPAPRPVVNVSGGPLLYSHRLSELRLLFGARDGAGSEHQINH 151

Query: 150

QAFSGEVQLIHYNHELYTNVTEAAKSPNGLVVVSIFIKVSDSSNPFLNRMLNRDTITRIT 209

QFS EVQLIH+N ELY N + A++ PNGL ++S+F+ V+

+SNPFL+R+LNRDTITRI+

55 Sbjct: 152 QGFSAEVQLIHFNQELYGNFSAASRGPNGLAILSLFVNVASTSNPFLSRLLNRDTITRIS 211

```
Query:
              570
     YKNDAYLLQGLNIEELYPETSSFITYDGSMTIPPCYETASWIIMNKPVYITRMQMHSLRL 2b9
                  YKNDAY LQ L++E L+PE+ FITY GS++ PPC ET +WI++++ + IT
     +QMHSLRL
     Sbjct:
              575
     YKNDAYFLQDLSLELLFPESFGFITYQGSLSTPPCSETVTWILIDRALNITSLQMHSLRL 271
              270 LSQNQPSQIFLSMSDNFRPVQPLNNRCIRTNINFSLQGKDC--PNNR 314
                  LSQN PSQIF S+S N RP+QPL +R +R N +
                                                        + C PN R
10
    Sbjct:
              272 LSQNPPSQIFQSLSGNSRPLQPLAHRALRGNRDPRHPERRCRGPNYR 31A
                 Pedant information for DKFZphamy2_15e14, frame 3
15
                          Report for DKFZphamy2_16e14-3
     ELENGTH
                     328
     EWWI
                     37563-19
20
     [pI]
                     8.22
     EHOMOLI
                     PIR:JE0375 carbonic anhydrase-related protein -
    human le-101
     EBFOCKZ
                     DMOLLO9B
    EBFOCK2
                     BL00162F
25
    EBFOCK2
                     Broozese
    EBFOCK21
                     BF007P5D
    EBFOCK21
                     BL00162C Eukaryotic-type carbonic anhydrases
    proteins
    EBFOCKZI
                     BL00162A Eukaryotic-type carbonic anhydrases
30
    proteins
    [[CCOP]
                     dlznca_ 2.56.1.1.3 Carbonic anhydrase [human
    (Homo sapiens le-103
    [SCOP]
                     d2cba___ 2.56.1.1.2 Carbonic anhydrase Ehuman
    (Homo sapiens 9e-97
35
    EEC]
                     4.2.1.1 Carbonate dehydratase le-36
    EEC B
                     3.1.3.48 Protein-tyrosine-phosphatase 2e-20
    EPIRKUI
                     blocked amino end &e-29
                     carbon-oxygen lyase le-36
    EPIRKUJ
    EPIRKUI
                     zinc le-36
40
    EPIRKUI
                     polymorphism 2e-20
    CPIRKWI
                     hydro-lyase le-36
    CPIRKWI
                     transmembrane protein 3e-23
    EPIRKUJ
                    tyrosine-specific phosphatase 2e-2D
    EPIRKUJ
                    brain be-16
45
    EPIRKWI
                     acetylated amino end le-36
    [PIRKW]
                    phosphatidylinositol linkage 2e-19 .
    EPIRKU
                    receptor 2e-20
                    liver 3e-29
    EPIRKWI
    EPIRKWI
                    phosphoprotein 2e-20
50
    [PIRKW]
                    saliva 2e-21
    CPIRKUJ
                    glycoprotein 2e-22
    EPIRKWI
                    mitochondrion le-32
    [PIRKW]
                    monomer 3e-32
    CPIRKWI
                    alternative splicing be-16
55
    EPIRKU
                    lipoprotein 2e-19
    EPIRKU
                    pyroglutamic acid 2e-21
    EPIRKU
                    metalloprotein Le-35
```

muscle 4e-31

EPIRKUI

WO 01/98454 PCT/IB01/02050 EPIRKWI. membrane protein 2e-19 CPIRKW3 phosphoric monoester hydrolase 2e-20 **EPIRKWI** homodimer 3e-23 fibronectin type III repeat homology 2e-20 ESUPFAMI 5 **ESUPFAMJ** carbonic anhydrase homology le-3b **ESUPFAMI** protein-tyrosine-phosphatase, receptor type zeta be-lb ESUPFAMI . carbonate dehydratase le-36 **ESUPFAMI** protein-tyrosine-phosphatase, receptor type gamma 10 2e-20 **ESUPFAMD** protein-tyrosine-phosphatase homology 2e-20 EZUPFAMI leukocyte common antigen cytosolic domain homology 2e-20 EPFAM3 Eukaryotic-type carbonic anhydrases 15 [KW] All_Beta **EKW3** ΔE SIGNAL_PEPTIDE 22 20 SEQ MEIVWEVLFLLQANFIVCISAQQNSPKIHEGWWAYKEVVQGSFVPVPSFWGLVNSAWNLC 25 **SVGKRQSPVNIETSHMIFDPFLTPLRINTGGRKVSGTMYNTGRHVSLRLDKEHLVNISGG** Jugc- ..TTTTCCCEETTTTTEETTTTCEEEEETT-TTCEEEEETTTTEEEEECTTTTTEEEEE 30 SFO PMTYSHRLEEIRLHFGSEDSQGSEHLLNGQAFSGEVQLIHYNHELYTNVTEAAKSPNGLV lugc- TTCCCEEEEEEEETTTTTTCTTTEETTBCCCEEEEEEEGG-GTTHHHHHHCTTTTEE 35 VVSIFIKVSDSSNPFLNRMLNRDTITRITYKNDAYŁLQGLNIEELYPETSSFITYDGSMT EEEEEEEC-CCCGGGHHHH--HHGGGCCTTTEEEETTTTCGGGGCCCCCCEEEEEECCC 40 SEQ IPPCYETASWIIMNKPVYITRM@MHSLRLLS@N@PS@IFLSMSDNFRPV@PLNNRCIRTN TTTTCCCEEEEECCCEEECHHHHHHHHHCCBCCTTTTCCCBTTTTCCCCCCTTTTCCEEC 45 INFSLQGKDCPNNRAQKLQYRVNEWLLK SFO lugc-(No Prosite data available for DKFZphamy2_1be14⋅3) 50 Pfam for DKFZphamy2_16e14-3 HMM_NAME Eukaryotic-type carbonic anhydrases 55

*WCYgeHWGPEHH....WHkhYPIAW....GDRQSPINIQWkearYDPS

PCT/IB01/02050

WO 01/98454

DKFZphamy2_lcl2

5 group: nucleic acid management

DKFZphamy2_1c12 encodes a novel 422 amino acid protein with partial identity to I-kappa-B-related protein and to BRCA1.

I-kappa-B-related protein interacts with transcription factors and BRCAL has a function in DNA damage response. I-kappa-B-alpha mutations contribute to constitutive NF-kappaB activity in cultured and primary HRS (Hodgkin/Reed-Sternberg) cells and are therefore involved in the pathogenesis of Hodgkin's disease (HD) patients.

The new protein can find application in modulating DNA repair and mutagenesis and also in expression profiling in HD related syndroms.

similarity to I-kappa-B-related protein

Sequenced by MediGenomix

25

30

20

Locus: unknown

Insert length: 1645 bp
Poly A stretch at pos. 1626, polyadenylation signal at pos. 1605

L GGATTTTCCT TGGTCTTAAG ATGGGTAGAA ATGTGATGCG ACACATGTCT

51 GATGACTTAG GAAGTTATGT TTCTCTTTCG TGTGATGACT TTTCTCACA

101 GGAATTAGAG ATTTTCATTT GCTCCTTTTC CTCCTCTGG CTTCAAATGT

35 L51 TTGTTGCAGA GGCAGTCTTT AAAAAAGTTGT GTCTACAGAG CTCTGGCAGT

201 GTTTCTTCTG AGCCACTCTC TCTTTCAGAAA ATGGTATATT CCTATTTACC

251 AGCCTTGGGG AAAACTGGTG TGCTTGGGTC TGGAAAGATT CAGGTGTCAA

301 AGAAAATAGG ACAGCGGCCT TGTTTTGACT CTCAGAGAAC CTTACTAATG

351 CTGAATGGTA CTAAAACAAAA ACAAGTCGAA GGGCTGCCAG AGTTACTAGA

40 401 CCTGAACCTT GCTAAATGTT CCTCATCATT AAAAAAAATTG AAAAAGAAGT

451 CAGAAGGAGA ATTGTCATGT TCCAAGGAGA ATTGCCCCTC TGTAGTTAAA

501 AAGATGAATT TTCACAAGAC TAATCTAAAA GGAGAAACAG CCTTGCATAG

551 AGCTTGCATA AATAACCAAG TAGAGAAATT GATTCTTCTT TCTCTTTTGC

451 CAGGAATAGA CATCAATGTT AAAGACAAAT CTGGCTGGAC GCCTTTGCAT

452 ATGATGCACT GTCAAACGGA CACAGTGTGT GTCCAGGAAA TTTTGCAACG

701 TTGTCCAGAG GTAGATCTGC TCACTCAAGT GGACGAGAC ACCCTTTTGC

751 ATGATGCACT GTCAAACGGA CATGAGAAA TTGGCCAGC GCTCTTTGC

851 CTTGGATTAT GTGGTTTCAC CTCAAACGAA TTGGCAAGCT GCTACTACAG

801 CAGAAAATAGA AGATACAGTG GAGAACTTC TTTGGCAACGT GCTACTACAG

801 CAGAAAATAGA AGATACAGTG GAGAACTTT ACAACAGAGG AATGCCTAAGG AGAAACAT

951 TTTCATTAT GTGGTTTCAC CTCAAACGAA TTGGCAAGCT GCTACTACAG

851 CTTGGATTAT TTGTTCAAATTT TTGATTTATC TTCAGAGTTC ATTTTACTTA

50 901 CAAAAATAGA AGATACAGTG GAGAACTTTC ATGCCACAGC AGAGAAACAT

951 TTTCATTACC AGCAACTTGA ATTTTGGCTTC TTTTACTTA GTAGGATGTT

1001 GCTAAATTTT TGTTCAATTT TTGATTTATC TTCAGAGTTC TTTAAGCTT

1001 AAAGAAACCA CCAGTGTTCA AATGAACTGC TTTTACTTA GTAGGATGCT

101 AAAGAAACCA CCAGTGTTCA AATGAACACA CTGAGGCCTT GTAAAGCTG

55 1151 AAAATTAAAAG ACATTGCACA AACTCCCACA CATTCTTAAG GAACTGCCTG

1201 AGAATTGAA AGCTTCCTAA ATTTGACTTCAAG TTTTAAGCTT

1201 AGAATTTGAA AGCTTCCTA AATTGACACA CTGAGGCCTT GATGATAACA

1201 AGAATTGAA AGCTTCCTA AATTGACACA CTGAGGCCTT GATGATAACA

1201 AGAATTGAA TGGTTCTCT GGGGTCACACA CTGAGGCCTT GATGATAACA

1201 AGAATTGAA TGGTTCTTCAA ATTTTCACTTACTGT

1201 AGAATTGAAT GACTTTCTAA ATCTGTTCAG TTTTCATTTACTGT

1351 GGACTTCATA GCTTACTGAC AGATAGTAAT TTGATTTATT TATTGACAGA
1401 CTTTGCAGCC TTGCTAAATT TTAAAAGCAT TTTTAAAAAA ACTTCTACAA
1451 AACTCTAGTA TGGGCTTCTG ACTTTTTCCA GGGTGTAGAA TTTGACTCAA
1501 AAGTAAAAAT AATTTTGTTT TAGTATATTC TACTTTCATT AATGTTTTTT
1551 TGTTCTGAAA GTGATATTAT ATTGTACATG TAAAATTAAT TTAAATATTT
1601 TTTCAAATAA AAATGTAATG TCCTGTAAAA AAAAAAAAA

BLAST Results

10

No BLAST result

15

Medline entries

No Medline entry

20

Peptide information for frame 3

25 ORF from 21 bp to 1286 bp; peptide length: 422 Category: similarity to known protein Classification: Cell signaling/communication

1 MGRNVMRHMS DDLGSYVSLS CDDFSSQELE IFICSFSSSW LQMFVAEAVF
30 51 KKLCLQSSGS VSSEPLSLQK MVYSYLPALG KTGVLGSGKI QVSKKIGQRP
101 CFDSQRTLLM LNGTKQKQVE GLPELLDLNL AKCSSSLKKL KKKSEGELSC
151 SKENCPSVVK KMNFHKTNLK GETALHRACI NNQVEKLILL LSLPGIDINV
201 KDNAGWTPLH EACNYGNTVC VQEILQRCPE VDLLTQVDGV TPLHDALSNG
251 HVEIGKLLQ HGGPVLLQQR NAKGELPLDY VVSPQIKEEL FAITKIEDTV
35 301 ENFHAQAEKH FHYQQLEFGS FLLSRMLLNF CSIFDLSSEF ILASKGLTHL
351 NELLMACKSH KETTSVHTDW LLDLYAGNIK TLQKLPHILK ELPENLKVCP

40

BLASTP hits .

No BLASTP hits available

401 GVHTEALMIT LEMMCRSVME FS

45

Alert BLASTP hits for DKFZphamy2_lcl2, frame 3

PIR:A56429 I-kappa-B-related protein - human, N = 1, Score = 242, P =

4-Le-la 50

TREMBLNEW:AF038042_1 gene: "BARD1"; product: "BRCA1-associated RING domain protein"; Homo sapiens BRCA1-associated RING domain protein

55 (BARDL) gener exons 10, 11 and complete cds., N = 1, Score = 236, P = 6.9e-17

>PIR:A56429 I-kappa-B-related protein - human Length = 481.

5 HSPs:

> Score = 242 (36.3 bits), Expect = 4.6e-18, P = 4.6e-18 Identities = 52/118 (44%), Positives = 71/118 (60%)

10 Query: 15b PSVVKKMNFHKTNLKGETALHRACINNQVEKLILLLSLPGIDINVKDNAGWTPLHEACNY 215 +++ N GET LHRACI Q+ ++ L+

GWTPLHEACNY

354 PGAAKGSKWNRRNDMGETLLHRACIEGQLRRVQDLVR-Sbict:

15 QGHPLNPRDYCGWTPLHEACNY 412

> 216 GNTVCVQEILQRCPEVDLL--Querv: TQVDGVTPLHDALSNGHVEIGKLLLQHGGPVLLQQRNA 272

G+ V+ +L VD +G+TPLHDAL+ GH E+ +LLL+ G V

20 L+ R A Sbjct: 413 GHLEIVRFLLDHGAAVDDPGGQGCEGITPLHDALNCGHFEVAELLLERGASVTLRTRKA 471

25 Pedant information for DKFZphamy2_lcl2, frame 3

Report for DKFZphamy2_1cl2.3

30 ELENGTHI 422 47071.18 EMWI 6.57 [[q] EHOMOLI PIR: A56429 I-kappa-B-related protein - human 3e-19 35 99 unclassified proteins ES. cerevisiae, YILLL2w1 **EFUNCATI** 3e-ll EFUNCATI 06.13.01 cytoplasmic degradation ES. cerevisiae. YGR232w3 4e-06 40 **IFUNCATI** 30.10 nuclear organization ES. cerevisiae, YIRD33w1 2e-04 EFUNCATD 04-05-01-07 chromatin modification ES. cerevisiae. YIR033w1 2e-04 **ESCOPI** dlawcb_ 1.91.3.1.2 GA binding protein (GABP) alpha 45 GA bindini be-24 EEC] 3.1.3.53 Myosin-light-chain-phosphatase 9e-Db **EPIRKWI** phosphotransferase 3e-07

EPIRKU tandem repeat 9e-06 **EPIRKU** transmembrane protein 7e-10 50 serine/threonine-specific protein kinase 3e-07 **EPIRKU**

EPIRKUI phosphoprotein 3e-10 **EPIRKWI** integrin binding 3e-07 **EPIRKWI** alternative splicing 3e-11 **EPIRKWI** peripheral membrane protein 2e-09

55 **EPIRKW** transcription regulation 3e-06 phosphoric monoester hydrolase 9e-06 **EPIRKUJ EPIRKWI** cvtoskeleton 4e-10

EPIRKU smooth muscle 9e-06

5	ESUPFAMD CERUD DE CER	ankyrin 3e-ll ankyrin repeat hom unassigned ankyrin Ank repeat [rregular]] _OW_COMPLEXITY	n repeat protei	ns 7e-1B	
10	SEQ MGRNVM SEG LawcB	1RHMSDDLGSYVSLSCD1	DFSS@ELEIFICSFS	SSSWL@MFVAEAVFKK	
15		SL@KMVYSYLPALGKT			· ·
20	SEG	DLNLAKCZZZLKKLKKI	· · · · · · · · · · · · · · · · · · ·	• • • • • • • • • • • • • • • • • • • •	
25	JawcB	CLILLLSLPGIDINVKD	• • • • • • • • • • • • • • • • • • • •	• • • • • • • • • • • • • • • • • • • •	•••••
30	JawcB	LZNGHVEIGKLLL@HG@			••••
35	SEQ ENFHAQ SEG LawcB	AEKHFHYQQLEFGSFLI	SRMLLNFCSIFDLS	· · · · · · · - · · · · · · ·	LLMACKSH
40	SEG KETTSV	HTDWLLDLYAGNIKTL	KLPHILKELPENLK	VCPGVHTEALMITLE	MMCRSVME
45	SEQ FS SEG lawcB -	·			·
50	(No Prosite	data available f	or DKFZphamy2_	_lcl2·3)	1
		Pfam	for DKFZphamy2	!_lcl2.3	
55		nk repeat			4.
	HMM	*GyTPLHIA	ARYNNVEMVrlLLQ	H-GADIN*	•

G+T+LH A+++N+VE LLL+ G DIN

Query 171 GETALHRACINNQVEKLILLLSLPGIDIN 199

34.48 (bits) f: 205 t: 232 Target: dkfzphamy2_lcl2.3
5 similarity to I-kappa-B-related protein
Alignment to HMM consensus:

Query *GyTPLHIAARyNNvEMVrlLLQHGADIN*
G+TPLH A+ Y+N+ +V+ LQ+ + ++

dkfzphamy2 205 GWTPLHEACNYGNTVCVQEILQRCPEVD 232

10

Query f: 239 t: 266 Target: dkfzphamy2_lcl2.3 similarity to I-kappa-B-related protein Alignment to HMM consensus:

#GyTPLHIAARyNNvEMVrlLLQHGADIN*

G TPLH A +++VE+ +LLLQHG +

Query 239 GVTPLHDALSNGHVEIGKLLLQHGGPVL 266

DKFZphamy2_lil

20

group: nucleic acid management

DKFZphamy2_lil encodes a novel L29 amino acidprotein with similarity to the murine hemin-sensitive initiation factor 2-

The hemin-sensitive initiation factor 2 is expressed predominantly in liver, spleen, colon and uterus and contains 2 protein kinase motifs. The mouse homologue inhibits protein synthesis in stress conditions by phosphorylation of eif-2-alphatour different eIF2alpha kinases have been identified in mammalian cells, the heme-regulated inhibitor (HRI), the interferon-inducible RNA-dependent kinase (PKR), the endoplasmic reticulum-resident kinase (PERK) and MGCN2. The new protein represents a new member of this family

The new protein can find application in modulating/blocking of translation.

40

similarity to hemin-sensitive initiation factor 2 (Mus musculus) $_{\text{\tiny T}}$ complete cds-alpha kinase

complete cds.

45 probably complete in genomic clone DJ0042M02

Sequenced by MediGenomix

Locus: /map="37.2 cR from top of Chr? linkage group"

50

Insert length: 2863 bp
Poly A stretch at pos. 2844, polyadenylation signal at pos. 2824

```
201 AAAAGAACCC CTACAACAGC CAACCTTCCC TTTTGCAGTT GCAAACCAAC
                                     TCTTGCTGGT TTCTTTGCTG GAGCACTTGA GCCACGTGCA TGAACCAAAC

CCACTTCGTT CAAGACAGGT GTTTAAGCTA CTTTGCCAGA CGTTTATCAA

ASL AATGGGGCTG CTGTCTTCTT TCACTTGTAG TGACGAGTTT AGCTCATTGA

GACTACATCA CAACAGAGCT ATTACTCACT TAATGAGGTC TGCTAAAGAG

ASL AGAGTTCGTC AGGATCCTTG TGAGGATATT TCTCGTTATCC AGAAAATCAG

ATCAAGGGAA GTAGCCTTGG AAGCACAAAC TTCACGTTAC TTAAATGAAT

TTGAAGAACT TGCCATCTTA GGAAAAAGGTG GATACGGAAG AGTATACAAG

BOL GTCAGGAATA AATTAGATGG TCAGTATTAT GCAATAAAAA AAATCCTGAT

ASL TAAGGGTGCA ACTAAAACAG TTTGCATGAA GGTCCTACGG GAAGTGAAGG

TTGCTGGCAGG TCTTCAGCAC CCCAATATTG TTGGCTATCA CACCGCGTGG

TGCTGGCAGG TCTTCAGCAC CCCAATATTG TTGGCTATCA CACCGCGTGG

TAAGAACATG TTCATGTGAT TCAGCCACGA GACAGAGCTG CCATTGAGTT

BOL GCCATCTCTG GAAGTGCTCT CCGACCAGGA AGAGGACAGA GAGCAATGTG

BSL GTGTTAAAAA TGATGAAAGT AGCAGCTCAT CCATTATCTT TGCTGAGCCC

TOL ACCCCAGAAA AAGAAAAACG CTTTGGAGAA TCTGACACTG AAAATCAGAA
                                         251 TCTTGCTGGT TTCTTTGCTG GAGCACTTGA GCCACGTGCA TGAACCAAAC
         5
    10
                              B5L GTGTTAAAAA TGATGAAAGT AGCAGCTCAT CCATTATCTT TGCTGAGCCC 90L ACCCCAGAAA AAGAAAAACG CTTTGGAGAA TCTGACACTG AAAATCAGAA 95L TAACAAGTCG GTGAAGTACA CCACCAATTT AGTCATAAGA GAATCTGGTG LOOL AACTTGAGTC GACCCTGGAG CTCCAGGAAA ATGGCTTGGC TGGTTTGTCT LO5L GCCAGTTCAA TTGTGGAACA GCAGCTGCCA CTCAGGCGTA ATTCCCACCT LLOL AGAGGAGAGT TTCACATCCA CCGAAGAATC TTCCGAAGAA AATGTCAACT LL5L TTTTGGGTCA GACAGAGGCA CAGTACCACC TGATGCTGCA CATCCAGATG L2OL CAGCTGTG AGCTCTCGCT GTGGGATTGG ATAGTCGAGA GAAACAAGCG L25L GGGCCGGGAG TATGTGGACG AGTCTGCCTG TCCTTATGTT ATGGCCAATG L3OL TTGCAACAAA AATTTTTCAA GAATTGGTAG AAGGTGTGTT TTACATACAT L35L AACATGGGAA TTGTGCACCG AGATCTGAAG CCAAGAAATA TTTTTCTTCA L4OL TGGCCCTGAT CAGCAAGTAA AAATAGGAGA CTTTGGTCTG GCCTGCACAG
   15
  20
                AACATGGGAA TTGTGCACCG AGATCTGAAG CCAAGAAATA TTTTTCTTCA
LHOL TGGCCCTGAT CAGCAAGTAA AAATAGGAGA CTTTGGTCTG GCCTGCACAG
LH5L ACATCCTACA GAAGAACACA GACTGGACCA ACAGAAACGG GAAGAGAACA
L5DL CCAACACATA CGTCCAGAGT GGGTACTTGT CTGTACGCTT CACCCGAACA
L55L GTTGGAAGGA TCTGAGTATG ATGCCAAGTC AGATATGTAC AGCTTGGGTG
LLOL TGGTCCTGCT AGAGCTCTTT CAGCCGTTTG GAACAGAAAT GGAGCGAGCA
LL5L GAAGTTCTAA CAGGTTTAAG AACTGGTCAG TTGCCGGAAT CCCTCCGTAA
L7DL AAGGTGTCCA GTGCAAGCCA AGTATATCCA GCACTTAACG AGAAGGAACT
L75L CATCGCAGAG ACCATCTGCC ATTCAGCTGC TGCAGAGTGA ACTTTTCCAA
LBDL AATTCTGGAA ATGTTAACCT CACCCTACAG ATGAAGATAA TAGAGCAAGA
LB5L AAAAGAAATT GCAGAACTAA AGAAGCAGCT AAACCTCCTT TCTCAAGACA
L9DL AAGGGGTGAG GGATGACGGA AAGGATGGGG GCGTGGGATG AAAGTGGACT
L75L TAACTTTTAA GGTAGTTAAC TGGAATGTAA ATTTTTAATC TTTATTAGGG
L9DL TATAGTTGGT ACAATGCTTC GTTGTATTTA GTAAGCCTTT ACAAGACTTG
L9DL TATAGTTGGT ACAATGCTCC CAAGCTGCCG TTCCTTCCCT TCCTGCCCCA
LDDL CAAGCTCCTT TTCCTGAATT TCCTACCTAA ATATTAACCA TATGCCTAGT
 25
 30
                         2101 CAAGCTCCTT TTCCTGAATT TCCTACCTAA ATATTAACCA TATGCCTAGT
40
45
50
                             2851 AAAAAAAAA AAA
```

55

Entry AFO288O8 from database EMBL:
Mus musculus hemin-sensitive initiation factor 2 alpha kinase
mRNA₁

5 complete cds.

Score = 6688, P = 2.7e-296, identities = 1922/2534

Entry ACOO5995 from database EMBL:
Homo sapiens clone DJOO42MO2, WORKING DRAFT SEQUENCE, 13
unordered
pieces.

Score = 5116, P = 0.0e+00, identities = 1090/1148

15

10

Medline entries

99042009:

20 Berlanga J.J., Herrero S., de Haro C.; Characterization of the hemin-sensitive eukaryotic initiation factor Zalpha kinase from mouse nonerythroid cells; J. Biol. Chem. 273(48):32340-32346(1998).

25

Peptide information for frame 1

30

ORF from 52 bp to 1938 bp; peptide length: 629 Category: similarity to known protein Classification: Protein management Prosite motifs: PROTEIN_KINASE_ATP (173-196) PROTEIN_KINASE_ST (437-449)

FOR KEIAELKKOL NLLZODKGVR DDGKDGGVG

1 MQGGNSGVRK REEEGDGAGA VAAPPAIDFP AEGPDPEYDE SDVPAEIQVL
40 51 KEPLQQPTFP FAVANQLLLV SLLEHLSHVH EPNPLRSRQV FKLLCQTFIK
101 MGLLSSFTCS DEFSSLRLHH NRAITHLMRS AKERVRQDPC EDISRIQKIR
151 SREVALEAQT SRYLNEFEEL AILGKGGYGR VYKVRNKLDG QYYAIKKILI
201 KGATKTVCMK VLREVKVLAG LQHPNIVGYH TAWIEHVHVI QPRDRAAIEL
251 PSLEVLSDQE EDREQCGVKN DESSSSSIIF AEPTPEKEKR FGESDTENQN
45 301 NKSVKYTTNL VIRESGELES TLELQENGLA GLSASSIVEQ QLPLRRNSHL
351 EESFTSTEES SEENVNFLGQ TEAQYHLMLH IQMQLCELSL WDWIVERNKR
401 GREYVDESAC PYVMANVATK IFQELVEGVF YIHNMGIVHR DLKPRNIFLH
451 GPDQQVKIGD FGLACTDILQ KNTDWTNRNG KRTPTHTSRV GTCLYASPEQ
501 LEGSEYDAKS DMYSLGVVLL ELFQPFGTEM ERAEVLTGLR TGQLPESLRK
50 551 RCPVQAKYIQ HLTRRNSSQR PSAIQLLQSE LFQNSGNVNL TLQMKIIEQE

55

BLASTP hits

No BLASTP hits available

trehalose)

Alert BLASTP hits for DKFZphamy2_lil, frame 1

No Alert BLASTP hits found

5 Pedant information for DKFZphamy2_lil, frame 1 ------

Report for DKFZphamy2_lil.1

10 ELENGTHD 646 EMWI 72738 - 78 [[q] 5-80 SWISSNEW: HRI_MOUSE HEME-REGULATED EUKARYOTIC []OMOL] INITIATION FACTOR EIF-2-ALPHA KINASE (EC 2.7.1.-) (HEME-REGULATED 15 INHIBITOR) (HRI) (HEME-CONTROLLED REPRESSOR) (HCR) (HEMIN-SENSITIVE INITIATION FACTOR-2 ALPHA KINASE). D.D. **EFUNCATE** 05.07 translational control ES. cerevisiae, YDR283c3 2e-43 20 ES. cerevisiae, YDR283cl 2e-43 EFUNCATI 10.02.11 key kinases ES. cerevisiae, YOR231w1 8e-14 [FUNCAT] 03.04 budding, cell polarity and filament formation ES. cerevisiae, YOR23Lwl &e-14 EFUNCATI D3.D1 cell growth ES. cerevisiae, YOR231w1 &e-14 ll.Ol stress response [S. cerevisiae, YOR231w] Be-14 **EFUNCATI EFUNCATI** Y0R231w1 8e-14 EFUNCATI 30.10 nuclear organization ES. cerevisiae, YKLlOlwl 30 8e-12 IFUNCATI 99 unclassified proteins ES. cerevisiae, YPL150wl 8e-15 EFUNCATI 03-13 meiosis ES- cerevisiae, YDR523cI 2e-11
EFUNCATI 03-10 sporulation and germination ES- cerevisiae, 35 YDR523cl 2e-ll 09.01 biogenesis of cell wall [FUNCAT] ES. cerevisiae. YPL14Ocl 4e-11 **EFUNCATI** 10.03.11 key kinases [S. cerevisiae, YCR073c] 9e-11 EFUNCATI 98 classification not yet clear-cut ES. cerevisiae. 40 YHRO82c3 le-lo [FUNCAT] 03.07 pheromone response, mating-type determination, sex-specific proteins ES. cerevisiae, YLR362w3 2e-10 **IFUNCATI** 10.99 other signal-transduction activities 45 cerevisiae, YDL101c1 3e-10 [FUNCAT] 11.04 dna repair (direct repair, base excision repair [FUNCAT] 04.05.01.01 general transcription activities 50 cerevisiae, YDLLO8wl le-09 EFUNCATI 03.16 dna synthesis and replication ES. cerevisiae, YBR160wl le-09 [FUNCAT] 01.05.04 regulation of carbohydrate utilization cerevisiae, YLR113w1 4e-09

ES. cerevisiae, YPLO31cl le-D8

- 5 [FUNCAT] c energy conversion [M- genitalium, MGl09] 2e-08 [FUNCAT] 03.19 recombination and dna repair [S. cerevisiae, Y0R35]c] le-07 [FUNCAT] 03.22.01 cell cycle check point proteins [S.
 - cerevisiae, YPL153cl le-0?

 CEFUNCATO 10.05.09 regulation of g-protein activity ES.

cerevisiae, YBLO16wl 7e-O7

EFUNCATU 04-O3-99 other trna-transcription activities

ES
cerevisiae, YILO35cl le-O6

EFUNCATI D8.13 vacuolar transport ES. cerevisiae, YGL18Dwl

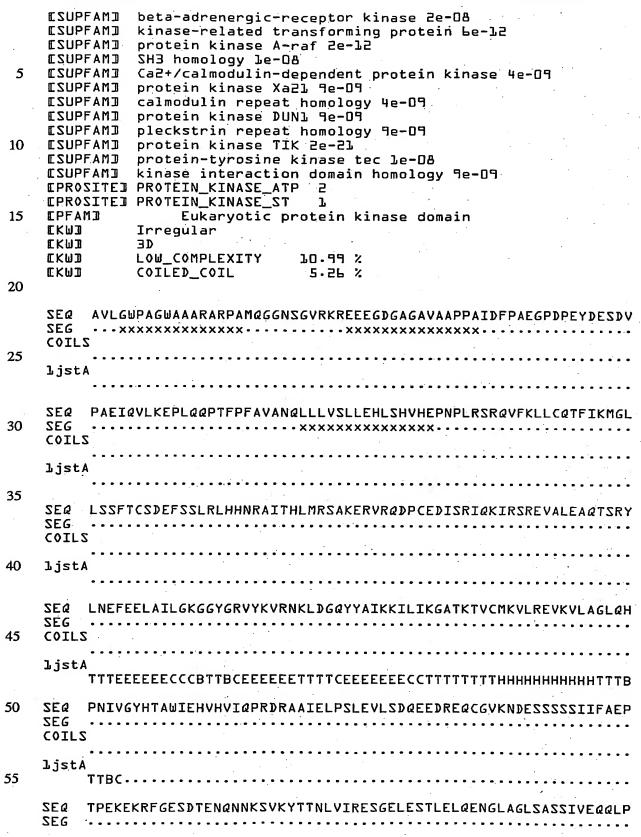
- 15 Le-Ob

 EFUNCATO Ob-13-O4 lysosomal and vacuolar degradation EScerevisiae, YGLIAOwo le-Ob

 EFUNCATO O4-99 other transcription activities ES- cerevisiae,
 YER129wo 2e-Ob
- 20 [FUNCAT] 30.02 organization of plasma membrane [S. cerevisiae, YDR122w] 2e-06 [FUNCAT] 30.07 organization of endoplasmatic reticulum [S. cerevisiae, YHR079c] 3e-06 [FUNCAT] 01.06.10 regulation of lipid, fatty-acid and sterol
- 25 biosynthesis ES. cerevisiae, YHRO79c3 3e-Ob EFUNCAT3 O8.99 other intracellular-transport activities ES. cerevisiae, YKL198c3 le-O5 EFUNCAT3 10.04.99 other nutritional-response activities ES. cerevisiae, YKL198c3 le-O5
- 30 IFUNCATI D9.04 biogenesis of cytoskeleton IS. cerevisiae.
 YNLO20cl 9e-05
 IFUNCATI Ob.07 protein modification (glycolsylation, acylation, myristylation, palmitylation, farnesylation and processing)
 IS. cerevisiae. YFLO33cl 4e-04
- 35 EFUNCATI 01.02.04 regulation of nitrogen and sulphur utilization
 ES. cerevisiae. YNL183cl 7e-04
 EBL0CKSl BL00107A Protein kinases ATP-binding region proteins
 ESCOPI dlir3a_ 5.1.1.2.6 insulin receptor Complex
 (transferase/substrate) le-22
- 40 [SCOP] dlfgkb_ 5.l.l.2.5 Fibroblast growth factor receptor l [human (Hom 9e-27] [SCOP] dlphk__ 5.l.l.l.6 gamma-subunit of glycogen phosphorylase kinas 2e-23 [SCOP] dlabo__ 5.l.l.l.4 Protein kiase (K2, alpha
- 45 subunit EMaize (Ze le-23
 ESCOPI d3lck__ 5-l-1-2-2 Lymphocyte kinase (lck) EHuman
 (Homo sapiens) 3e-22
 ESCOPI d2erk__ 5-l-1-1-11 MAP kinase Erk2 Erat (Rattus norvegicus) 7e-20
- 50 [SCOP] dlcdkb_ 5.l.l.l.2 cAMP-dependent PK, catalytic subunit Comple Le-19 [SCOP] dlhcl_ 5.l.l.l. Cyclin-dependent PK [Human (Homo sapiens) 5e-2]
 - TECD 2.7.1.112 Protein-tyrosine kinase Le-O8
- 55 [EC] 2.7.1.126 beta-Adrenergic-receptor kinase 2e-D8 [EC] 2.7.1.117 Myosin-light-chain kinase 1e-D9
 - EECl 2.7.1.37 Protein kinase 5e-12

```
[EC]
               2.7.1.123 Ca2+/calmodulin-dependent protein kinase 4e-
    09
     EPIRKWI .
                    phosphotransferase 0.0
    CPIRKUD
                   nucleus 9e-09
    EPIRKU
                    RNA binding 2e-21
     EPIRKWI
                    duplication &e-10
     [PIRKW]
                    tandem repeat 4e-09
     [PIRKW]
                    zinc 5e-12
    EPIRKWJ
                    cell cycle control 2e-09
10
    IPIRKWI
                    serine/threonine-specific protein kinase D.O
    EPIRKWI.
                    transmembrane protein 2e-09
                    zinc finger 8e-10
    CPIRKWI
    [PIRKW]
                    oncogene be-12
    EPIRKWD
                    autophosphorylation [] []
15
    EPIRKU
                    coat protein le-11
    EPIRKUI
                   magnesium 9e-09
    EPIRKUJ
                    ATP 0-0
    EPIRKUJ
                    polyprotein be-12
    EPIRKU
                   receptor 9e-09
20
    EPIRKUJ
                   phosphoprotein D.D
    EPIRKUJ
                   sporulation 2e-09
    EPIRKWI
                   glycoprotein 9e-09
    EPIRKU
                   growth factor receptor 9e-11
    EPIRKU
                   signal transduction 2e-12
25
    EPIRKWI
                   serine/threonine/tyrosine-specific protein kinase
    8e-10
    EPIRKU
                   protein kinase 8e-10
    EPIRKUI
                   transforming protein 2e-12
    EPIRKWI
                   heme binding D.O
30
    EPIRKUJ
                   purine nucleotide binding 2e-10
    EPIRKWI
                   calcium binding 4e-09
    EPIRKUJ
                   meiosis le-O8
                   alternative splicing le-ll
    [PIRKW]
                   P-loop Ze-10
    EPIRKU
35
    LPIRKUJ
                   proto-oncogene 2e-12
    EPIRKUI
                   segmentation 4e-10
    EPIRKWI
                   stress-induced protein le-09
    EPIRKU
                   EF hand 4e-09
    EPIRKUI
                   cell division le-09
40
    EPIRKUD
                   calmodulin binding 4e-09
    EZUPFAMI
              LIM protein kinase &e-10
    EZUPFAMI
              calcium-dependent protein kinase 4e-89
    ESUPFAMD
              rat protein kinase raf 5e-12
              AMP-activated protein kinase 2e-08
    ESUPFAMI
45
    ESUPFAMD
              protein kinase byr2 5e-09
    ESUPFAMD
              SH2 homology le-OA
              unassigned Ser/Thr or Tyr-specific protein kinases 0.0
    EZUPFAMJ
              leucine-rich alpha-2-glycoprotein repeat homology 9e-09
    ESUPFAMI
50
              double-stranded RNA-binding repeat homology 2e-21
    ESUPFAMI
              histidine--tRNA ligase homology be-42
    EZUPFAMI
    ESUPFAMI
              SAM homology 5e-09
    ESUPFAMD
              avian retrovirus ICLO gag-Rmil-env polyprotein le-ll
    ESUPFAMI
              LIM metal-binding repeat homology &e-10
55
    ESUPFAMI
              GCN2 protein Le-42
    ESUPFAMJ
              protein kinase homology [].[]
              protein kinase C zinc-binding repeat homology 2e-12
    ESUPFAMI
    ESUPFAMI
              Ca2+/calmodulin-dependent protein kinase II 4e-D8
```





нмм

10

30

35

HMM

Query

233 KVLAGLQHPNIVGYHTAWI-EHVHV 256

*IYMIMEYMeGGDLFDYIrrng······pMsEweIrfIMyQIL +++ M+++E +L+D+I++++

+++

Query 396 LHIQMQLCELSLWDWIVERNKRGREYVDAZ3AVATKIFQELV 443

15 HMM
rGMeYLHSMgIIHRDLKPENILIDeN-gqIKIcDFGLARqMn.....
+G+ Y+H+MGI+HRDLKP+NI++ + Q+KI+DFGLA+
Query 444
EGVFYIHNMGIVHRDLKPRNIFLHGPDQQVKIGDFGLACTDILQKNTDWT 493

20
HMM
----nYerMttfCGTPWYMMAPEVIImgnyYttkVDMWSFGCILWEMMT
+ T+++GT Y +PE ++G++Y+ K+DM+S+G++L

E++
25 Query 494 NRNGKRTPTHTSRVGTCLYA-SPEQ-LEGSEYDAKSDMYSLGVVLLELF- 540

GepPFyd.-dnMemImrIiqr-frrpfWpnCSeElyDFMrwCWnyDPekR +PF ++ E + ++ ++ ++ +C+ +++ + + + ++R

Query 541 -- QPFGTEMERAEVLTGLRTGQLPESLRKRCPVQAKYIQ-HLTRRNSSQR 587

P++ Q+L++ F
Query 588 PSAIQLLQSELF 599

PTFr@ILnHPWF*

DKFZphamy2_lil4

5 group: transmembrane proteins

DKFZphamy2_lil4 encodes a novel 617 amino acid protein with similarity to the human 1(3)mbt protein homolog.

- 10 Mutations of the Drosophila 1(3)mbt gene lead to malignant brain tumors. The novel protein contains 1 transmembrane domain.

 No informative BLAST results: No predictive prosite: pfam or SCOP motife
- The new protein can find application in studying the expression profile of oncogenes and amgydala-specific genes and as a new marker for amygdala cells.
- 20 similarity to Human 1(3)mbt protein homolog mRNA
 - > 14 exons (HS756G23 (EMBLNEW))
 Pedant: TRANSMEMBRANE 1
- 25 Sequenced by MediGenomix

Locus: /map="22ql3.3l-l3.33"

Insert length: 3071 bp

30 Poly A stretch at pos. 3052, no polyadenylation signal found

1 GGCAGGCCAA TATGGCTTCC TGCACCTGGT GACGCTTGGC GAAACTGAGG
51 TCTCATGGAG AAGCCCCGGA GTATTGAGGA GACCCCATCT TCAGAACCAA
35 101 TGGAGGAAGA GGAAGATGAC GACTTGGAGC TGTTTGGTGG CTATGATAGT
151 TTCCGGAGTT ATAACAGCAG TGTTGGGCAGT GAGAGCAGGT CCTATCTTGAA
201 GGAGTCAAGT GAAGCAGAAA ATGAGGATCG GGAAGCAGGG GAACTGCCGA
251 CCTCCCCGCT GCATTTGCTC AGCCCTGGGA CTCCTCGCTC CTTGGATGGC
301 AGTGGTTCTG AGCCAGCTGT CTGTGAGATG TGTGGTATCG TGGGTACAAG
40 351 GGAAGCCTTC TTCTCCAAGA CCAAAGAGGT CTGCCGACTC TCCTGCTC
401 GGAGCTACTC CTCCAACTCC AAGAAAAGCCA GTATCTTGGC TAGAATTACAG
451 GGAAAACCAC CGACCAAAAA AGCCAAAGGT CTGCACAGG CTGCCTGGTC
501 TGCCAAAAAT GGAGCCTTCC TCCACTCCA AGGACAGGA CAGCTGGCAG
551 ATGGGACACC AACAAGAAA GACCCACTCTG TCTTGGATGAC
45 L01 AAGTTCCTGA AGGACCAAG TTACAAGGCT GCTCCCGTC GCTGGTC
501 TGCGAGGTCC AACAAGGACAA GACGCTCTGG TCTTGGGGGG
45 L01 AAGTTCCTGA AGGACCAAG TTACAAGGCT GCTCCCGTC GCTGTTTCAA
L51 GCCTCTGTCA TCCAGACAG TTACAAGGCT GCTCCCGTC GCTGTTTCAA
L51 GCCTCTGTCA TCCAGACAG AGGGAAGAGA TGTGCTGGATG
50 A51 ATGCCACC CATTGGCTG TGTGCCTC CCAGCCGGGT GTACTGGATC
50 A51 ATGCCACCC CATTGGCTG TGTGCCACCA GCAACCTG GGAACAGGA
51 ACGGCTGGTG GGCTCCAGGA CGCTTCCCGT GGAACCTG GGAACAGTG
51 ACGGCTGGTG GGCTCCAGGA CGCTTCCCGT GGAACCTG GGAACAGTG
52 ACGGCTGGT GGCTCCAGGA CGCTTCCCGT GGAACCTG GGAACAGTG
53 ACGCCAGCC ATCCATGCAA GTTCACCGAC TGGAAGGGT ACCTCATGAA
54 ACGCCAGCC TTTAGACCA GTTCACCGAC GGAACAGGT GACACAGAT
55 ATGCCACCC CATTGGCAG CCCTTCCCGT GGAATTCCAC
56 ACAAGTCCC AGGACTACCC TTTAGGCAGG GCATGCGGCT GGAAGTGGTG
57 ACGGCTGGT GGCTCCAGGA CCCCTTCCCGT GGAATTCCAC
57 ACGGCTGGT GGCTCCAGGA CCCCTGCATC GCATTCCCGT GGAAGGGT ACCACAGATAT
55 ACGCCTGGCC CTACGGCTC TCTACGGAGA TGGGGCT GCAAGGACGACT
1151 TCTGGTGCCA CACGGCTCC TCTACGGAGA TGGGGCT GCCACGGCT
1151 TCACCCCCAC TTCCGGAAGGA TCCACTGGAGG ACCACGGACT
1251 TCACCCCCAC TTCCGGAAGA TCCACTGTCCT TACCTCTCA
1251 TCACCCCCAC TTCCCGGAAGA TCCACTCTCCT TACCCCCAC TTCCCGTTCCT TACCTCTCA

WO 01/98454 PCT/IB01/02050 1301 AGAAGGTACG AGCAGTCTAC ACAGAAGGCG GTTGGTTTGA GGAAGGGATG 1351 AAGCTGGAGG CCATTGACCC CCTGAATCTG GGCAACATCT GCGTGGCAAC 1401 TGTCTGTAAG GTTCTCCTGG ATGGATACCT GATGATCTGT GTGGACGGGG 1451 GGCCCTCCAC AGATGGCTTG GACTGGTTCT GCTACCATGC CTCTTCCCAC 5 1501 GCCATCTTCC CGGCCACCTT CTGTCAGAAG AATGACATTG AGCTCACACC 1551 GCCAAAAGGT TATGAGGCAC AGACTTTCAA CTGGGAGAAC TACTTGGAGA LLUL AGACCAAGTC GAAAGCCGCT CCATCGAGAC TCTTTAACAT GGATTGCCCA
LLUL AACCATGGCT TCAAGGTGGG CATGAAGCTG GAGGCCGTGG ACCTGATGGA 10 15 20 25 30 2753 CCTCTCTGTG TAAATTCTGC CCGGTGCTGT GAAGGCTGGA CGGTGGAGGA 2801 CCTGCTGGGG TCTCCTGGGA CCCGCCTGTT GCTTCTGCCC TCCCCTGTGG 2851 AAAGGTCTAT ATGACGGGCC GCCTGAGGCC CCAGAACTCG TCTGTGAACC 2901 ACCTTTCCA GCCAGAGTTC CCAAAGCTGG AACGCTAGCT GCCTGCTCTT 2951 CCTTAAGATG GCCTCCCCC GACCCGCCAC GGCCCTCAGT TGCCAGGGAT

BLAST Results

40

35

Entry HS756623 from database EMBLNEW: Human DNA sequence from clone 756623 on chromosome 22ql3.3l-l3.33 Score = 3939_7 P = $0.0e+00_7$ identities = 875/954

BDD1 GGGGCCACCA CTGTCACACT GTGGAATACA AGACAGTGAA CTCTGTCTGC

45

Entry UB9358_1 from database TREMBL: product: "1(3)mbt protein homolog"; Human 1(3)mbt protein homolog mRNA; complete cds.

50 Score = 505, P = 7.2e-45, identities = 123/320, positives = 170/320,

3051 CTAAAAAAA AAAAAAAAA A

frame +1

Entry ABO14581_1 from database TREMBL:
55 gene: "KIAAO681"; product: "KIAAO681 protein"; Homo sapiens
mRNA for
KIAAO681 protein; partial cds.

Score = 503_1 P = 1.4e-4b, identities = 122/307, positives = 163/307, frame +1

5

Medline entries

10 No Medline entry

Peptide information for frame 1

ORF from 55 bp to 1905 bp; peptide length: 617 Category: similarity to known protein Classification: unclassified

20

25

30

MEKPRSIEET PSSEPMEEEE DDDLELFGGY DSFRSYNSSV GSESSSYLEE

51 SSEAENEDRE AGELPTSPLH LLSPGTPRSL DGSGSEPAVC EMCGIVGTRE

101 AFFSKTKRFC SVSCSRSYSS NSKKASILAR LQGKPPTKKA KVLHKAAWSA

151 KIGAFLHSQG TGQLADGTPT GQDALVLGFD WGKFLKDHSY KAAPVSCFKH

201 VPLYDQWEDV MKGMKVEVLN SDAVLPSRVY WIASVIQTAG YRVLLRYEGF

251 ENDASHDFWC NLGTVDVHPI GWCAINSKIL VPPRTIHAKF TDWKGYLMKR

301 LVGSRTLPVD FHIKMVESMK YPFRQGMRLE VVDKSQVSRT RMAVVDTVIG

351 GRLRLLYEDG DSDDDFWCHM WSPLIHPVGW SRRVGHGIKM SERRSDMAHH

401 PTFRKIYCDA VPYLFKKVRA VYTEGGWFEE GMKLEAIDPL NLGNICVATV

451 CKVLLDGYLM ICVDGGPSTD GLDWFCYHAS SHAIFPATFC QKNDIELTPP

501 KGYEAQTFNW ENYLEKTKSK AAPSRLFNMD CPNHGFKVGM KLEAVDLMEP

551 RLICVATVKR VVHRLLSIHF DGWDSEYDQW VDCESPDIYP VGWCELTGYQ

601 LQPPVAAGVG SRGPKRL

35

BLASTP hits

No BLASTP hits available

40

Alert BLASTP hits for DKFZphamy2_lil4, frame l

TREMBL:ABD14581_1 gene: "KIAAD681"; product: "KIAAD681 protein"; Homo

- 45 sapiens mRNA for KIAAOb&l protein, partial cds., N = l, Score = 503, P = 3.9e-48
- TREMBL:UA9358_1 product: "1(3)mbt protein homolog": Human 1(3)mbt protein homolog mRNA; complete cds.; N = 1; Score = 505; P = 6.2e-48
- 55 >TREMBL:UB9358_1 product: "1(3)mbt protein homolog": Human 1(3)mbt protein homolog mRNA, complete cds.

 Length = 772

HSPs:

Score = 505 (75.8 bits): Expect = 6.2e-48: P = 6.2e-48 5 Identities = 123/313 (39%), Positives = 170/313 (54%)

293 WKGYLMKRLVGSRTLPVDFH--IKMVESMKYPFRQGMRLEVVDKSQVSRTRMAVVDTVIG 35D

W+ YL ++ + T PV + V K F+ GM+LE +D

V V G 10 Sbjct: 509 MEZATEEGK--AITAPVSLFQDSQAVTHNKNGFKLGMKLEGIDPQHPSMYFILTVAEVCG 265

351 GRLRLLYEDGDSD-DDFWCHMWSPLIHPVGWSRRVGHGIKMSE--RRSDMAHHPTFRKIY 407

15 RLRL + DG S+ DFW + SP IHP GW + GH +++ + 266 YRLRLHF-Sbjct:

DGYSECHDFWNNASPDIHPAGWFEKTGHKLQLPKGYKEEEFSWSQYMCSTR 324

20 Query: 4D8 CDAVP-YLFKKVRAVYTEGGWFEEGMKLEAIDPLNLGNICVATVCKVLLDGYLMICVDGG 466 G F+ GMKLEA+D +N ++ D

Sbjct: 325 AQAAPKHMFVSQSHSPPPLG-FQVGMKLEAVDRMNPSLVCVASVTDVV-25 DSRFLVHFDNW 382

467 PSTDGLDWFCYHASSHAIFPATFCQKNDIELTPPKGY-Query: EAGTFNWENYLEKTKSKAAPSR 525

D++CZZ I P +CQK LTPP+ Y + F WE YLE+T

30 Sbict: 383 DDT--YDYWC-DPSSPYIHPVGWCQKQGKPLTPPQDYPDPDNFCWEKYLEETGASAVPTW 439

Query:

40

35 LFNMDCPNHGFKVGMKLEAVDLMEPRLICVATVKRVVHRLLSIHFDGWDSEYDQWVDCES 585 F + P H F V MKLEAVD P LI VA+V+ V YD W+D +

440 AFKVR-Sbjct: PPHSFLVNMKLEAVDRRNPALIRVASVEDVEDHRIKIHFDGWSHGYDFWIDADH 49A

586 PDIYPVGWCELTGYQLQPPV 605 Querv: PDI+P GWC TG+ LQPP+ 499 PDIHPAGWCSKTGHPLQPPL 518 Sbjct:

Score = 333 (50.0 bits), Expect = 4.le-27, P = 4.le-27 45 Identities = 103/324 (31%), Positives = 151/324 (46%)

179 FDWGKFLKDHSYKAAPVSCFKHVPLYDQWEDVMK-GMKVEVLNZDAVLPSRVYWIASVIQ 237

50 + 🛭 +L++ APVS F+ ++ K GMK+E + D PS . +Y+I +V + Sbjct: 20b WSWESYLEEQKAITAPVSLFQDSQAVTHNKNGFKLGMKLEGI--DPQHPS-MYFILTVAE 262

55 Query: 23B TAGYRVLLRYEGFENDASHDFWCNLGTVDVHPIGWCAINSKILVPPRTIHAKFTDWKGYL 297 HDFW N + D+HP GW GYR+ L ++G+ Y+

Sbjct: 2b3 VCGYRLRLHFDGYSE-CHDFWYDAHIDGUAHIDGYSAHIDGYSAHVWHDDH

Query: 298 MKRLVGSRTLPVDFHIKMVESMKYP---

5 FRQGMRLEVVDKSQVSRTRMAVVDTVIGGRLR 354

+R H+ +S P F+ GM+LE VD+ S +A V

V+ R

Sbjct: 321 CS----

TRÄQAAPKHMFVSQSHSPPPLGFQVGMKLEAVDRMNPSLVCVASVTDVVDSRFL 37b

10

20

Query: 355 LLYEDGDSDDDFWCHMWSPLIHPVGWSRRVGHGIKMSERRSD---MAHHPTFRKIYCDAV 411

+ +++ D D+WC SP IHPVGW ++ G + + D

AV

15 Sbjct: 377

VHFDNWDDTYDYWCQPZSPYIHPVGWCQKQGKPLTPPQDYPDPDNFCWEKYLEETGASAV 436

Query: 412

PYLFKKVRAVYTEGGWFEEGMKLEAIDPLNLGNICVATVCKVLLDGYLMICVDGGPSTDG 471

P KVR ++ F MKLEA+D N I VA+V V D + I

DG + G Sbjct: 437 PTWAFKVRPPHS----FLVNMKLEAVDRRNPALIRVASVEDVE-DHRIKIHFDGW--SHG 489

25 Query: 472 LDWFCYHASSHAIFPATFCQKNDIELTPPKG 502
D F A I PA +C K L PP G
Sbjct: 490 YD-FWIDADHPDIHPAGWCSKTGHPLQPPLG 519

Score = 236 (35.4 bits), Expect = 2.5e-16, P = 2.5e-16 30 Identities = 47/110 (42%), Positives = 66/110 (60%)

Query: 499 PPKGYEAQTFNWENYLEKTKSKAAPSRLF-NMDCPNH--- GFKVGMKLEAVDLMEPRLIC 554

P G + + ++WE+YLE+ K+ AP LF + H GFK+GMKLE +D

35 P +

Sbjct: 197

PATGEKKECUSWESYLEEQKAITAPVSLFQDSQAVTHNKNGFKLGMKLEGIDPQHPSMYF 256

Query: 555 VATVKRVVHRLLSIHFDGWDSEYDQWVDCESPDIYPVGWCELTGYQLQPP

40 604

+ TV V L +HFDG+ +D WV+ SPDI+P GW E TG++LQ P Sbjct: 257 ILTVAEVCGYRLRLHFDGYSECHDFWVNANSPDIHPAGWFEKTGHKLQLP

30P

45

Pedant information for DKFZphamy2_lil4, frame l

Report for DKFZphamy2_lil4.1

50

ELENGTHI 617

EMM3 69264-33

EpII 6-05

55 EHOMOLI TREMBL:UB9358_1 product: "1(3)mbt protein homolog"; Human 1(3)mbt protein homolog mRNA; complete cds. le-47

EBLOCKSI BLO1206A Amiloride-sensitive sodium channels proteins

	[KW]	TRANSMEMBRANE 1 LOW_COMPLEXITY 9-40 %
5	SEQ SEG PRD MEM	MEKPRSIEETPSSEPMEEEEDDDLELFGGYDSFRSYNSSVGSESSSYLEESSEAENEDRExxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
	SEQ SEG PRD MEM	AGELPTSPLHLLSPGTPRSLDGSGSEPAVCEMCGIVGTREAFFSKTKRFCSVSCSRSYSS
15	SEQ SEG PRD MEM	NSKKASILARL@GKPPTKKAKVLHKAAWSAKIGAFLHS@GTG@LADGTPTG@DALVLGFDxxxxxcchhhhhhhhhhhhhhhhhhhhhhhhhhhh
20	SEQ SEG PRD MEM	WGKFLKDHSYKAAPVSCFKHVPLYDQWEDVMKGMKVEVLNSDAVLPSRVYWIASVIQTAG
25	SEQ SEG PRD MEM	YRVLLRYEGFENDASHDFWCNLGTVDVHPIGWCAINSKILVPPRTIHAKFTDWKGYLMKR
30	SEQ SEG PRD MEM	LVGSRTLPVDFHIKMVESMKYPFRQGMRLEVVDKSQVSRTRMAVVDTVIGGRLRLLYEDG
35	SEQ SEG PRD MEM	DSDDDFUCHMUSPLIHPVGUSRRVGHGIKMSERRSDMAHHPTFRKIYCDAVPYLFKKVRA
40	SEQ SEG PRD MEM	VYTEGGWFEEGMKLEAIDPLNLGNICVATVCKVLLDGYLMICVDGGPSTDGLDWFCYHAS cccccchhhhhheeeecccccceeeeeeehhhhhceeeeee
45	SEQ SEG PRD MEM	SHAIFPATFCQKNDIELTPPKGYEAQTFNWENYLEKTKSKAAPSRLFNMDCPNHGFKVGM CCCCCCCCCCCCCCCCCChhhhhhhhhhhhhhhhhhh
50	SEQ SEG PRD MEM	KLEAVDLMEPRLICVATVKRVVHRLLSIHFDGWDSEYDQWVDCESPDIYPVGWCELTGYQ
55	SEQ SEG PRD MEM	L@PPVAAGVGSRGPKRL

(No Prosite data available for DKFZphamy2_lil4.1)

5 (No Pfam data available for DKFZphamy2_lil4.l)

DKFZphamy2_li24

group: differentiation/development

DKFZphamy2_li24 encodes a novel 835 amino acid protein without partial similarity to rattus norvegicus Notch2 protein.

10 Notch family molecules are thought to be negative regulators of neuronal differentiation in early brain development. Notch2 is expressed not only by neuronal cells in the embryonic brain, but also by glial cells in the postnatal brain. The new protein represents a new member of this family and may be involved in specific differentation or developmental pathways of the nervous system.

The new protein can find application in modulating development and differentiation of amygdala cells.

20

putative protein

probably complete cds.

25

35

Sequenced by MediGenomix

Locus: unknown

30 Insert length: 2768 bp Poly A stretch at pos. 2714, polyadenylation signal at pos. 2697

L AGAAATCTTC AGCCAAACAG CTGCAGGAAG TAGAGAAGGT TAAACCCCAG

51 AGTGAGAAAG TTCATCAGAC TCTGATTCTG GACCCAGCAC AGAGGAAGAG LOL ACTCCAGCAG CAGATGCAGC AGCACGTTCA GCTCTTGACC CAAATCCACC LSL TTCTTGCCAC CTGCAACCCC AACCTCAATC CGGAGGCCAC TACCACCAGG 201 ATATTTCTTA AAGAGCTGGG AACCTTTGCT CAAAGCTCCA TCGCCCTTCA 251 CCATCAGTAC AACCCCAAGT TTCAGACCCT GTTCCAACCC TGTAACTTGA 40 BOL TGGGAGCTAT GCAGCTGATT GAAGACTTCA GCACACATGT CAGCATTGAC 351 TGCAGCCCTC ATAAAACTGT CAAGAAGACT GCGAATGAAT TTCCCTGTTT 401 GCCAAAGCAA GTGGCTTGGA TTCTGGCCAC AAGCAAGGTT TTCATGTATC 451 CAGAGTTACT TCCAGTGTGT TCCCTGAAGG CAAAGAATCC CCAGGATAAG 501 ATCGTCTCA CCAAGGCTGA GGACAATTTG TTAGCTTTAG GACTGAAGCA 551 TTTTGAAGGA ACTGAGTTTC CTAATCCTCT AATCAGCAAG TACCTTCTAA 601 CCTGCAAAAC TGCCCACCAA CTGACAGTGA GAATCAAGAA CCTCAACATG 45 LSD AACAGAGCTC CTGACAACAT CATTAAATTT TATAAGAAGA CCAAACAGCT 7D1 GCCAGTCCTA GGAAAATGCT GTGAAGAGAT CCAGCCACAT CAGTGGAAGC 751 CACCTATAGA GAGAGAAGAA CACCGGCTCC CATTCTGGTT AAAGGCCAGT 50 BDL CTGCCATCCA TCCAGGAAGA ACTGCGGCAC ATGGCTGATG GTGCTAGAGA B51 GGTAGGAAAT ATGACTGGAA CCACTGAGAT CAACTCAGAT CGAAGCCTAG 901 AAAAAGACAA TTTGGAGTTG GGGAGTGAAT CTCGGTACCC ACTGCTATTG 951 CCTAAGGGTG TAGTCCTGAA ACTGAAGCCA GTTGCCACCC GTTTCCCCAG IDDI GAAGGCTTGG AGACAGAAGC GTTCATCAGT CCTGAAGCCC CTCCTTATCC 55 1051 AACCCAGCCC CTCTCTCCAG CCCAGCTTCA ACCCTGGGAA AACACCAGCC LLDL CGATCAACTC ATTCAGAAGC CCCTCCGAGC AAAATGGTGC TCCGGATTCC 1151 TCACCCAATA CAGCCAGCCA CTGTTTTACA GACAGTTCCA GGTGTCCCTC 1201 CACTGGGGGT CAGTGGAGGT GAGAGTTTTG AGTCTCCTGC AGCACTGCCT

· WO 01/98454 PCT/IB01/02050 1251 GCTGTGCCCC CTGAGGCCAG GACAAGCTTC CCTCTGTCTG AGTCCCAGAC 1301 TTTGCTCTCT TCTGCCCCTG TGCCCAAGGT AATGCTGCCC TCCCTTGCCC 1351 CTTCTAAGTT TCGAAAGCCA TATGTGAGAC GGAGACCCTC AAAGAGAAGA 1401 GGAGTCAAGG CCTCTCCCTG TATGAAACCT GCCCCTGTTA TCCACCACCC 1451 TGCATCTGTT ATCTTCACTG TTCCTGCTAC CACTGTGAAG ATTGTGAGCC 1501 TTGGCGGTGG CTGTAACATG ATCCAGCCTG TCAATGCGGC TGTGGCCCAG 5 1551 AGTCCCCAGA CTATTCCCAT CACTACCCTC TTGGTTAACC CTACTTCCTT 1601 CCCCTGTCCA TTGAACCAGT CCCTTGTGGC CTCCTCTGTC TCACCCTTAA 1651 TTGTTTCTGG CAATTCTGTG AATCTTCCTA TACCATCCAC CCCTGAAGAT 1701 AAGGCCCACG TGAATGTGGA CATTGCTTGT GCTGTGGCTG ATGGGGAAAA 10 1751 TGCCTTTCAG GGCCTAGAAC CCAAATTAGA GCCCCAGGAA CTATCTCCTC 1801 TCTCTGCTAC TGTTTTCCCG AAAGTGGAAC ATAGCCCAGG GCCTCCACTA 1851 GCAGATGCAG AGTGCCAAGA AGGATTGTCA GAGAATAGTG CCTGTCGCTG LODI GACCETTETE AAAACAGAGG AGGGGAGGCA AGCTCTGGAG CCGCTCCCTC 15 1951 AGGGCATCCA GGAGTCTCTA AACAACCCTA CCCCTGGGGA TTTAGAGGAA 2DD1 ATTGTCAAGA TGGAACCTGA AGAAGCTAGA GAGGAAATCA GTGGATCCCC 2051 TGAGCGTGAT ATTTGTGATG ACATCAAAGT GGAACATGCT GTGGAATTGG 2101 ACACTGGTGC CCCAAGCGAG GAGTTGAGCA GTGCTGGAGA AGTAACGAAA . 2151 CAGACAGTCT TACAGAAGGA AGAGGAGAGG AGTCAGCCAA CTAAAACCCC 2201 TTCATCTTCT CAAGAGCCCC CTGATGAAGG AACCTCAGGG ACAGATGTGA 20 2251 ACAAAGGATC ATCAAAGAAT GCTTTGTCCT CAATGGATCC TGAAGTGAGG 2301 CTTAGTAGCC CCCCAGGGAA GCCAGAAGAT TCATCCAGTG TTGATGGTCA 2351 GTCAGTGGGG ACTCCAGTTG GGCCAGAAAC TGGAGGAGAG AAGAATGGGC 2401 CAGAAGAAGA GGAAGAAGAG GACTTTGATG ACCTCACCCA AGATGAGGAA 2451 GATGAAATGT CATCAGCTTC TGAGGAATCT GTGCTTTCTG TCCCAGAACT 25 2501 CCAGGTGAGA GCTGGAGAAT ATTCTCAAGT ATTTCGTGGA CTCAGTAATA 2551 TGTATCACTT ATTGATATGC CACCTGCTTG CTTGCTGCAC TATGGATAGT 2601 CCTAAAATCA TTTGTATTTG ATTTGTGAAT GCATTATGGG ACATGATTGT 2651 GGAGTTGAGG TGAAATGAGA TGGAAAGGAT GAAATTTTAC TTATTATATT 30 2751 AAAAAAAA AAAAAAA

BLAST Results

35

Entry RNNOTCHX from database EMBL: Rat notch 2 mRNA.

Score = 818, P = 1.6e-26, identities = 216/277

40

Medline entries

45

No Medline entry

50

55

Peptide information for frame 3

ORF from 114 bp to 2618 bp; peptide length: 835 Category: putative protein Classification: Differentiation/Development

1 MQQHVQLLTQ IHLLATCNPN LNPEATTTRI FLKELGTFAQ SSIALHHQYN 51 PKFQTLFQPC NLMGAMQLIE DFSTHVSIDC SPHKTVKKTA NEFPCLPKQV

	WO 01/98454	PCT/IB01/02050
5	151 EFPNPLISKY LLTCKTAHQL 201 KCCEEIQPHQ WKPPIEREEH 251 TGTTEINSDR SLEKDNLELG	LKAKNPQDKI VFTKAEDNLL ALGLKHFEGT TVRIKNLNMN RAPDNIIKFY KKTKQLPVLG RLPFWLKASL PSIQEELRHM ADGAREVGNM SESRYPLLLP KGVVLKLKPV ATRFPRKAWR SFNPGKTPAR STHSEAPPSK MVLRIPHPIQ
	351 PATVLQTVPG VPPLGVSGGE	SFESPAALPA VPPEARTSFP LSESQTLLSS
	451 FTVPATTVKI VSLGGGCNMI	VRRRPSKRRG VKASPCMKPA PVIHHPASVI QPVNAAVAQS PQTIPITTLL VNPTSFPCPL
10	501 NQSLVASSVS PLIVSGNSVN	LPIPSTPEDK AHVNVDIACA VADGENAFQG VEHSPGPPLA DAECQEGLSE NSACRWTVVK
	LOI TEEGRAALEP LPAGIAESLN	NPTPGDLEEI VKMEPEEARE EISGSPERDI LSSAGEVTK@ TVL@KEEERS @PTKTPSSS@
•	701 EPPDEGTSGT DVNKGSSKNA	LSSMDPEVRL SSPPGKPEDS SSVDGQSVGT
15	BOI GEYSQVFRGL SNMYHLLICH	FDDLTQDEED EMSSASEESV LSVPELQVRA LLACCTMDSP KIICI
		BLASTP hits
20	No DIASTR bite	
73	No BLASTP hits available	
	Alert BLASTP h	its for DKFZphamy2_li24, frame 3
25	No Alert BLASTP hits found	
	Pedant informat	ion for DKFZphamy2_li24, frame 3
30	Report	for DKFZphamy2_li24.3
	ELENGTHI 872 EMWI 95366.29	
35	[pI] 5.87	waan 1 D alaka alaa isaa isaa
	3.2.1.3) - yeast (Saccharom)	ucan l-4-alpha-glucosidase (EC /ces cerevisiae) 5e-06
	<pre>EFUNCATI 30.01 organization YIR019cl 2e-07</pre>	
40	EFUNCATI 30.90 extracellula YIRO19cl 2e-07	ar/secretion proteins ES. cerevisiae.
		ate utilization [[S. cerevisiae]
45	EFUNCATI 02.10 tricarboxyli YDR148cI 5e-04	c-acid pathway ES. cerevisiae,
43	EFUNCATI 30.16 mitochondria	ıl organization [[S. cerevisiae]
	YDR148cl 5e-04 EKWl Alpha_Beta	
50	EKWI LOW_COMPLEXITY	9.40 %
	SEA VSELVALABLEVIVEABENING	
	ZEG	LILDPAQRKRLQQQMQQHVQLLTQIHLLATCNPNLNP
55	PRD ccchhhhhhhhhhhhccccchhhh	hcccchhhhhhhhhhhhhhhhhhhhhhhccccccc
	SEQ EATTTRIFLKELGTFAQSSIALH	HQYNPKFQTLFQPCNLMGAMQLIEDFSTHVSIDCSPH
		ccccceeeecccchhhhhhhhhhceeeeecccc

-	SEQ SEG PRD	eeeeeccccccchhhhhhhccceeeecccccccccccc
5	SEQ SEG PRD	LKHFEGTEFPNPLISKYLLTCKTAHQLTVRIKNLNMNRAPDNIIKFYKKTKQLPVLGKCC hheeecccccccceeeeeeeeeeeeeeehhhhhhhhheeecccccc
10	SEQ SEG PRD	EEIQPHQWKPPIEREEHRLPFWLKASLPSIQEELRHMADGAREVGNMTGTTEINSDRSLE eeecccccccchhhhhhcceeeecchhhhhhhhhhhhh
15	SEQ SEG PRD	KDNLELGSESRYPLLLPKGVVLKLKPVATRFPRKAWR@KRSSVLKPLLI@PSPSL@PSFNxxxxxxxxxxxxxxxxxxxxxxxxxx
20	SEQ SEG PRD	PGKTPARSTHSEAPPSKMVLRIPHPI@PATVL@TVPGVPPLGVSGGESFESPAALPAVPP
25	SEQ SEG PRD	EARTSFPLSES@TLLSSAPVPKVMLPSLAPSKFRKPYVRRRPSKRRGVKASPCMKPAPVI
	SEQ SEG PRD	HHPASVIFTVPATTVKIVSLGGGCNMIQPVNAAVAQSPQTIPITTLLVNPTSFPCPLNQS
30	SEQ SEG PRD	CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
35	SEQ SEG PRD	SATVFPKVEHSPGPPLADAEC@EGLSENSACRUTVVKTEEGR@ALEPLP@GI@ESLNNPT
40	SEQ SEG PRD	PGDLEEIVKMEPEEAREEISGSPERDICDDIKVEHAVELDTGAPSEELSSAGEVTKQTVL
45	SEQ SEG PRD	@KEEERS@PTKTPSSS@EPPDEGTSGTDVNKGSSKNALSSMDPEVRLSSPPGKPEDSSSV
	SEQ SEG PRD	DGQSVGTPVGPETGGEKNGPEEEEEEDFDDLTQDEEDEMSSASEESVLSVPELQVRAGEY
50	SEQ SEG PRD	
55		Prosite data available for DKFZphamy2_li24.3)
	(Na	Diam data available for DVF7nhamv2 li24.3)

5

DKFZphamy2_ljl9

group: differentiation/development

- 10 DKFZphamy2_ljl9 encodes a novel 15D amino acid protein with high similarity to the allograft inflammatory factor-1 of Cyprinus carpio.
- Allograft inflammatory factor-1 (AIF-19 is a protein involved in allograft rejection. In experimental autoimmune encephalomyelitis (EAE), neuritis(EAN) and uveitis (EAU) it is produced by macrophages and microglia cells.
- The new protein can find clinical application in the development of tools to enhance the compatibility of transplanted tissues as well as in expression profiling of autoimmune diseases and infections.
- 25 strong similarity to allograft inflammatory factor-1 (Cyprinus carpio)

identical to DKFZphamy2_lnl

30 Sequenced by MediGenomix

Locus: /map="504.9 cR from top of Chr9 linkage group"

Insert length: 3381 bp

35 Poly A stretch at pos. 3362, polyadenylation signal at pos. 3344

1 GCCGGAGCCC GGACCAGGCG CCTGTGCCTC CTCCTCGTCC CTCGCCGCGT 51 CCGCGAAGCC TGGAGCCGGC GGGAGCCCCG CGCTCGCCAT GTCGGGCGAG 40 LOL CTCAGCAACA GGTTCCAAGG AGGGAAGGCG TTCGGCTTGC TCAAAGCCCG 151 GCAGGAGAGG AGGCTGGCCG AGATCAACCG GGAGTTTCTG TGTGACCAGA 201 AGTACAGTGA TGAAGAGAAC CTTCCAGAAA AGCTCACAGC CTTCAAAGAG 251 AAGTACATGG AGTTTGACCT GAACAATGAA GGCGAGATTG ACCTGATGTC BOL TTTAAAGAGG ATGATGGAGA AGCTTGGTGT CCCCAAGACC CACCTGGAGA 351 TGAAGAAGAT GATCTCAGAG GTGACAGGAG GGGTCAGTGA CACTATATCC 45 401 TACCGAGACT TTGTGAACAT GATGCTGGGG AAACGGTCGG CTGTCCTCAA 451 GTTAGTCATG ATGTTTGAAG GAAAAGCCAA CGAGAGCAGC CCCAAGCCAG 501 TTGGCCCCC TCCAGAGAGA GACATTGCTA GCCTGCCCTG AGGACCCCGC 551 CTGGACTCCC CAGCCTTCCC ACCCCATACC TCCCTCCGA TCTTGCTGCC 50 LOL CTTCTTGACA CACTGTGATC TCTCTCTCTC TCATTTGTTT GGTCATTGAG 651 GGTTTGTTTG TGTTTTCATC AATGTCTTTG TAAAGCACAA ATTATCTGCC 701 TTAAAGGGC TCTGGGTCGG GGAATCCTGA GCCTTGGGTC CCCTCCTCT 751 CTTCTTCCCT CCTTCCCCGC TCCCTGTGCA GAAGGGCTGA TATCAAACCA BOL AAAACTAGAG GGGGCAGGGC CAGGGCAGGG AGGCTTCCAG CCTGTGTTCC B5L CCTCACTTGG AGGAACCAGC ACTCTCCATC CTTTCAGAAA GTCTCCAAGC POL CAAGTTCAGG CTCACTGACC TGGCTCTGAC GAGGACCCCA GGCCACTCTG 55 951 AGAAGACCTT GGAGTAGGGA CAAGGCTGCA GGGCCTCTTT CGGGTTTCCT LOOL TEGACAGTEC CATEGTTCCA GTECTCTEGT GTCACCCAGE ACACAGCCAC

```
1051 TCGGGGCCCC GCTGCCCCAG CTGATCCCCA CTCATTCCAC ACCTCTTCTC
         LIDI ATCCTCAGTG ATGTGAAGGT GGGAAGGAAA GGAGCTTGGC ATTGGGAGCC
         1151 CTTCAAGAAG GTACCAGAAG GAACCCTCCA GTCCTGCTCT CTGGCCACAC
         LEDI CTGTGCAGGC AGCTGAGAGG CAGCGTGCAG CCCTACTGTC CCTTACTGGG
         1251 GCAGCAGAGG GCTTCGGAGG CAGAAGTGAG GCCTGGGGTT TGGGGGGAAA
  5
         ACCOMPAGA STOATAGE SEATHTTO ACCTTTTAGE GAGACTA ACTGENTAGE GAGACTA ACCTTTTAGE GAGACTA ACCTTTTAGE GAGACTAGE ACCTTTTAGE ACCTTTTAGE GAGACTAGE ACCTTTTAGE ACCTTTAGE ACCTTTAGE
         1351 GATGGGAGAA TGAGGAGTAA AATGCTCACG GCAAAGTCAG CAGCACTGGT
         1401 AAGCCAAGAC TGAGAAATAC AAGGTTGCTT GTCTGACCCC AATCTGCTTG
         1501 CACTCATTGA CTCACTCATT CACCAGATAT TTATTGACCT GCTATTATAA
10
         1551 GCTTTACATC CTCCCATGTT GTCCTGGCAT GTGCAGTATA CACGGTCTAA
         1603 CTCATCTCTC CCCAGATCTC TCAGAACCTT GAGCTTGGGA ATTGAACTGG
         1651 GGTCACCTGT GTCCTTTCTT ATGGACTCGC AGGATTTTAG AACCCTAATG
         1701 CACCCTGGAG GGTAGCTGGG CCAGACTTCT CATTTCACAG GTGAGGAGAC
         1751 TGGTGCCCCA CAGGGATTAA GTGCCTTGCC CAAGGTCAGG CTTATCTCCA
15
         LBDL GAGGGAGGTG CCCTGGACTG GGGCCCAGAT GTTCAGGGAC CCTGCCTACA
         LBSL CCTCATTTCC AGTGTGGGCT GCCTTAGTTA GTTATGAGAA CAGGGAAGGG
         1901 CTGGGAAGAG ACAGCCTCCA AGGTCAACAC TTGGAGAGGG TTTCACTTGC
         1951 TCTGAAGACC CTGGTCCAGG ATTCGCCCTC TCCCATGCCT
                                                                                                  TCAAGTCAGC
         2001 ATCAGGCTTA GGGCAAAGAC CAGGCCTCTG AAGCTGCCTC TTGTAATTCA
2051 TGCAGGAAGA TGTCAAAGTC AGCCCCATCT TGGCTGATCA GGGTGTTCAG
2101 CCTTAACCCC ACCTGTGTTC TGAAGTCTCT TACCCTACCT GCTCAGGACT
2151 GAGACAGTTA TTCACTGAAC ATATTTATTA AGCACTTGCT GTAGGCCAAC
20
         2201 AGTTAAGAAT CCAATAATGA AATGGACAGA TTCATGGAAC
                                                                                                   TTAGAGTCCA
         2251 ATAGGAAAGT GAGACCCAGA CAATGACAAT GAGATAAATG TTAGGAAGGG
25
         2301 GGAGGTATGG GGTGACTTCC CTGCAGTCCT GGGGGCCTAC ATGGGCCCAA
2351 GACTGGGTGA GAGTCTTGGC AGAGCCTTTG CAACACCTTA AGTGGACAGG
         24D1 ACTGGGAGGT CTTGGTGGTT GGAGCCAACG TGGGTTCCCT GCGGCTCCTT
         2453 AGTCACCTCT GATAGCAGAT TGAGGGAGGA AAACAGGTAA GGCATGAGGA
30
         2501 AATGGCCAGG TTGGGTTAAC CCACTGGTTT CAACCAGTTC AGGAATGAGG
         2551 TTATTTGGCC ATGACTGGCT GATCTTGAGC TCAAGGATCT GCTTCAAATG
         2603 CACACAGGCC TAGTTGAAGT TTAAACCCCA GCAAAACATT CCTCCCTGTA
         2653 AATGGAAAAT CCTACTTCTA CCCCCACCCT GCCCTGTTTT TTGTTTTTTT
         2751 CTGGGACAGG CTGGGACCTT TGAGGAAGAT AAAGCCTTCC TTGACTACCC
35
         2801 ATCATATTCA GTGTCCCTGT TCCTCACTCA GAGAGGAAGG CAGAACCAGT
         2851 CAGGCTTATT TCAGTAAGTT CCACAGTTCT ACAAGACTGC AGGAATTCTC
         2901 CTTAAGGGAG GAGAGCAAGC AGGTGTGGCC CCAGCTTCTG GAAATGGCAG
         2951 AAGAGAGGGT TTTCTCATTG AATGGGGGTG GGGGCTCGTG TGTCCTGGGA
40
         BDD AACCCCATCA GTCCCTTCAT TTCTTGAGAC TCAACTCCTG GGAGAGAGG
         3051 GTCTCAAGAG TTGTCCCTGG AAGGAGGGCG GGGGCAGTCT GCATCTATTT
         3101 CAGGTTGTGG CTCTTGGTTC TAGGACTCTT ACTTCTCTGG CTAAGGGCTC
         3151 AGCTTCTTGG GACTTCAACC ATCTTCTTTC TGAAAGACCA AATCTAATGT
         ACCAGTAAC GTGAGGACTG CCAAGTATGG CTTTGTCCCT ATGACTCAGA
         3251 GGAGGGTTTG TCGGGCAAAT TCAGGTGGAT GAAGTATGTG TGTGCGTGTG
45
         3301 CATGGGAGTG TGCGTGGACT GGGATATCAT CTCTACAGCC TGCAAATAAA
         A AAAAAAAA AAAAAAAA A
```

50

BLAST Results

Entry ABO12309_1 from database TREMBL:
product: "allograft inflammatory factor-l": Cyprinus carpio mRNA

55 for allograft inflammatory factor-L, complete cds.

Score = 575, P = 3.7e-54, identities = 113/146, positives = 128/146,

frame +2

5

Medline entries

No Medline entry

10

Peptide information for frame 2

15 ORF from 89 bp to 538 bp; peptide length: 150 Category: strong similarity to known protein Classification: unclassified

1 MSGELSNRFØ GGKAFGLLKA RØERRLAEIN REFLCDØKYS DEENLPEKLT
20 51 AFKEKYMEFD LNNEGEIDLM SLKRMMEKLG VPKTHLEMKK MISEVTGGVS
101 DTISYRDFVN MMLGKRSAVL KLVMMFEGKA NESSPKPVØP PPERDIASLP

. 25

BLASTP hits

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_ljl9, frame 2

No Alert BLASTP hits found

Pedant information for DKFZphamy2_ljl9, frame 2

35

30

Report for DKFZphamy2_ljl9.2

ELENGTHI 150

40 EMWI 17067.86

EPII 6.63

EHOMOLI TRE

factor-1"; Cyprinu

factor-1 complete

45 EFUNCATI 30.04 ord

EHOMOLI TREMBL:ABOL2309_l product: "allograft inflammatory factor-l"; Cyprinus carpio mRNA for allograft inflammatory factor-l, complete cds. 2e-59

cerevisiae, YBRLD9cl 5e-04

EFUNCATI D3.22 cell cycle control and mitosis ES. cerevisiae, YBRLD9cl 5e-04

EFUNCATI 03.04 budding, cell polarity and filament formation 55

ES. cerevisiae, YBRL09cl 5e-04

EFUNCATI 03.01 cell growth ES. cerevisiae, YBRL09cl 5e-04

EFUNCATI 30.05 organization of centrosome ES. cerevisiae, YBRL09cl 5e-04

PCT/IB01/02050 WO 01/98454 **ESCOPI** d2mysb_ 1.37.1.5.15 Myosin Essential Chain Myosin Regulatory Chai 5e-20 **EZCOPI** dlwdcb_ 1.37.1.5.14 Myosin Essential Chain Myosin Regulatory Chai 3e-05 5 dlosa__ 1.37.1.5.13 Calmodulin E(Paramecium **EZCOPI** tetraurelia) 3e-lb ESCOPI dlauib_ 1.37.1.5.19 Calcineurin regulatory subunit (B-chain 2e-16 · [PIRKW] duplication 7e-06 10 [PIRKW] mitosis 7e-06 [PIRKW] calcium binding 7e-06 [PIRKW] EF hand 7e-06 **EPIRKU**J cell division 7e-06 **ESUPFAM3** unassigned calmodulin-related proteins 3e-47 calmodulin 7e-06 15 **ESUPFAMD** calmodulin repeat homology 3e-47 **EZUPFAM3** EKWI All_Alpha [KW] αE 20 SEQ MSGELSNRFQGGKAFGLLKARQERRLAEINREFLCDQKYSDEENLPEKLTAFKEKYMEFD lctrтнининининининин 25 SEQ LNNEGEIDLMSLKRMMEKLGVPKTHLEMKKMISEVTGGVSDTISYRDFVNMMLGKRSAVL lctr-ТТТТТСВСНИНИННИНТТТСССИНИННИННИНСТТТТСССВСИНИННИНССТТТИНИ SEQ KLVMMFEGKANESSPKPVGPPPERDIASLP 30 lctr-HHHHHHTTTTC...... (No Prosite data available for DKFZphamy2_ljl9.2) 35 (No Pfam data available for DKFZphamy2_ljl9.2)

DKFZphamy2_24b4

5 group: cell cyle

> DKFZphamy2_2464 encodes a novel 698 amino acid protein with similarity to human STIML.

- 10 The stromal interaction molecular 1 gene (STIM1) encodes a type I trans-membrane protein of unknown function, which induces growth arrest and degeneration of the human tumor cell lines 6401 and RD but not HBL100 and CaLu-6, suggesting a role in the pathogenesis of rhabdomyosarcomas and rhabdoid tumors. There is also strong similarity to a Mus musculus stromal cell protein, which
- selectively increases interleukin 7-dependent proliferation of pre-B cells. The novel protein contains 1 transmembrane domain.
- The new protein can find application in modulation of tumour .20 growth.

similarity to STIML (Homo sapiens)

25 probably differential polyadenylation: cf. EST-BLAST file. perhaps complete cds. Pedant: SIGNAL_PEPTIDE and TRANSMEMBRANE

Sequenced by GBF

30

35

Locus: /map="139.2 cR from top of Chr4 linkage group"

Insert length: 3305 bp Poly A stretch at pos. 3274, polyadenylation signal at pos. 3260

1 GGCGCCTTCA TCCCGCCTCG ACTCCTGGCC CAGCGTGGGG CTGGCTGCTG 51 CGGCGGCGGC GCTGGGCTGC GTTGCTGGTG CTCGGGCTGC TGGTACCCGG 101 AGCGGCGGAC GGATGCGAGC TTGTGCCCCG GCACCTCCGC GGGCGGCGGG 40 151 CGACTGGCTC TGCCGCAACT GCCGCCTCCT CTCCCGCCGC GGCGGCCGGC 201 GATAGCCCGG CGCTCATGAC AGATCCCTGC ATGTCACTGA GTCCACCATG 251 CTTTACAGAA GAAGACAGAT TTAGTCTGGA AGCTCTTCAA ACAATACATA **BOL AACAAATGGA TGATGACAAA GATGGTGGAA TTGAAGTAGA GGAAAGTGAT** 351 GAATTCATCA GAGAAGATAT GAAATATAAA GATGCTACTA ATAAACACAG 45 4D1 CCATCTGCAC AGAGAAGATA AACATATAAC GATTGAGGAT TTATGGAAAC 451 GATGGAAAAC ATCAGAAGTT CATAATTGGA CCCTTGAAGA CACTCTTCAG 501 TGGTTGATAG AGTTTGTTGA ACTACCCCAA TATGAGAAGA ATTTTAGAGA 551 CAACAATGTC AAAGGAACGA CACTTCCCAG GATAGCAGTG CACGAACCTT LOI CATTTATGAT CTCCCAGTTG AAAATCAGTG ACCGGAGTCA CAGACAAAAA 50 **L51** CTTCAGCTCA AGGCATTGGA TGTGGTTTTG TTTGGACCTC TAACACGCCC 7DL ACCTCATAAC TGGATGAAAG ATTTTATCCT CACAGTTTCT ATAGTAATTG 751 GTGTTGGAGG CTGCTGGTTT GCTTATACGC AGAATAAGAC ATCAAAAGAA **BOL CATGTTGCAA AAATGATGAA AGATTTAGAG AGCTTACAAA CTGCAGAGCA** B51 AAGTCTAATG GACTTACAAG AGAGGCTTGA AAAGGCACAG GAAGAAAACA PDL GAAATGTTGC TGTAGAAAAG CAAAATTTAG AGCGCAAAAT GATGGATGAA 55 951 ATCAATTATG CAAAGGAGGA GGCTTGTCGG CTGAGAGAGC TAAGGGAGGG LODL AGCTGAATGT GAATTGAGTA GACGTCAGTA TGCAGAACAG GAATTGGAAC 1051 AGGTTCGCAT GGCTCTGAAA AAGGCCGAAA AAGAATTTGA ACTGAGAAGC

LLDL AGTTGGTCTG TTCCAGATGC ACTTCAGAAA TGGCTTCAGT TAACACATGA 1151 AGTAGAAGTG CAATACTACA ATATTAAAAG ACAAAACGCT GAAATGCAGC 1201 TAGCTATTGC TAAAGATGAG GCAGAAAAA TTAAAAAGAA GAGAAGCACA 1251 GTCTTTGGGA CTCTGCACGT TGCACAGGC TCCTCCCTAG ATGAGGTAGA LOOP CCACAAAATT CTGGAAGCAA AGAAAGCTCT CTCTGAGTTG ACAACTTGTT 5 1351 TACGAGAACG ACTTTTTCGC TGGCAACAAA TTGAGAAGAT CTGTGGCTTT BUDD CAGATAGCCC ATAACTCAGG ACTCCCCAGC CTGACCTCTT CCCTTTATTC 1451 TGATCACAGC TGGGTGGTGA TGCCCAGAGT CTCCATTCCA CCCTATCCAA 1501 TTGCTGGAGG AGTTGATGAC TTAGATGAAG ACACACCCC AATAGTGTCA
1551 CAATTTCCCG GGACCATGGC TAAACCTCCT GGATCATTAG CCAGAAGCAG 10 JUDI CAGCCTGTGC CGTTCACGCC GCAGCATTGT GCCGTCCTCG CCTCAGCCTC JL5J AGCGAGCTCA GCTTGCTCCA CACGCCCCCC ACCCGTCACA CCCTCGGCAC 1701 CCTCACCACC CGCAACACAC ACCACACTCC TTGCCTTCCC CTGATCCAGA 1751 TATCCTCTCA GTGTCAAGTT GCCCTGCGCT TTATCGAAAT GAAGAGGAGG 1801 AAGAGGCCAT TTACTTCTCT GCTGAAAAGC AATGGGAAGT GCCAGACACA 15 1851 GCTTCAGAAT GTGACTCCTT AAATTCTTCC ATTGGAAGGA AACAGTCTCC
1901 TCCTTTAAGC CTCGAGATAT ACCAAACATT ATCTCCGCGA AAGATATCAA
1951 GAGATGAGGT GTCCCTAGAG GATTCCTCCC GAGGGGATTC GCCTGTAACT 2001 GTGGATGTGT CTTGGGGTTC TCCCGACTGT GTAGGTCTGA CAGAAACTAA 2051 GAGTATGATC TTCAGTCCTG CAAGCAAAGT GTACAATGGC ATTTTGGAGA 20 2101 AATCCTGTAG CATGAACCAG CTTTCCAGTG GCATCCCGGT GCCTAAACCT 2151 CGCCACACAT CATGTTCCTC AGCTGGCAAC GACAGTAAAC CAGTTCAGGA 2201 AGCCCCAAGT GTTGCCAGAA TAAGCAGCAT CCCACATGAC CTTTGTCATA 2251 ATGGAGAGAA AAGCAAAAAG CCATCAAAAA TCAAAAGCCT TTTTAAGAAG 25 2301 AAATCTAAGT GAACTGGCTG ACTTGATGGA ATCATGTTCA AGTGGCATCT 2351 GTAAACTATT ATCCCCCACC CTCCACTCCC CACCTTTTTT TTGGTTTAAT 2401 TTTAGGAATG TAACTCCATT GGGGCTTTCC AGGCCGGATG CCATAGTGGA 2451 ACATCCAGAA GGGCAACTGT CTACTGTCTG CTTATTTAAG TGACTATATA 2501 TAATCAATTC ATCAAGCCAG TTATTACTGA AAAATCATTG AAATGAGACA 2551 GTTTACAGTC ATTTCTGCCT ATTTATTTCT GCTTTGTTCT CAGTGATGTA 30 2601 TATGCAACAT TTTGTTGAAA GCCACGATGG ACTTACAAGC TTTAATGGAC 2651 TCGTAAGCCA GCATGGGCTT GCAAAAATTT CTTGTTTACC AGAGCATCTT 2701 CTTATCTTTC CACAGAGCTA TTTACATCCT GGACTATATA ACTTAAAAGA 2751 AGTAAAACGT AATTGCACTA CTGTTTTCCA GACTGGAAAA AAAAAAAAAT 35 2801 CTCTGCAAGT GAAACTGTAT AGAGTTTATA AAATGACTAT GGATAGGGGA 2851 CTGTTTTCAC TTTTAGATCA AAATGGGTTT TTAAGTAGAA CCTAGGGTTT 2901 CTAATTGACT TGATTTCTGG AAATGAAAAC CCGCGCTTTT ATTATGGGAA 2951 GCTTCTTGAA CTGCATTTAC TATTGTGAAG TTTCAAGTCC CGCTGTAAAG 3001 ATCATGTTGT TTTGTTTTCC CCAGGGCTTT CACTGTGATT TACTGCATTG 3051 CAGGCTGTAT GATAAAACAC ACATAATTTA AAGAGAGAAG GCTCTTGATT 40 BLOL CCTTATGCAA GTGGAAGAGT TGAAACTTGA TTGAAGGACT TAAAACATTC 3151 ACAACCTTAA GCCGAGGTGG GGGGATATGG GGATTCAGGC AGTTGTTTAC 3201 ACACTTTGAA TAACTGCAAA GGATTTACGG TTTGTGAAAA ATGTGTACTG 3251 TGGAAAGAT AATAAATTGA AGACATTAAA AAAAGAAAA AAAAAAAAA

BLAST Results

Entry HS5242610_1 from database TREMBL:
gene: "STIM1"; product: "GOK"; Homo sapiens GOK (STIM1) mRNA;
complete
cds.
Score = 1397; P = 4.2e-142; identities = 275/447; positives =

45

50

AAAA LOEE

Entry MMU47323_1 from database TREMBL:

product: "stromal cell protein": Mus musculus stromal cell
protein

mRNA, complete cds.

5 Score = 1394, P = 8.8e-142, identities = 274/447, positives = 336/447, frame +3

Entry HS917349 from database EMBL:

10 human STS EST167479.

Score = 1390, P = 9.1e-57, identities = 284/287

15

Medline entries

97079692:

Parker NJ, Begley CG, Smith PJ, Fox RM, Molecular cloning of a 20 novel

human gene (D1154896E) at

chromosomal region llpl5.5. Genomics 1996 Oct 15:37(2):253-6

96326680:

25 Oritani K, Kincade PW.; Identification of stromal cell products that interact with pre-B cells. J Cell Biol 1996 Aug;134(3):771-82

30

Peptide information for frame 3

35

ORF from 216 bp to 2309 bp; peptide length: 698 Category: similarity to known protein Classification: Cell signaling/communication Prosite motifs: RGD (589-591)

40

1 MTDPCMSLSP PCFTEEDRFS LEALQTIHKQ MDDDKDGGIE VEESDEFIRE 51 DMKYKDATNK HSHLHREDKH ITIEDLWKRW KTSEVHNWTL EDTLQWLIEF 101 VELPQYEKNF RDNNVKGTTL PRIAVHEPSF MISQLKISDR SHRQKLQLKA 151 LDVVLFGPLT RPPHNWMKDF ILTVSIVIGV GGCWFAYTQN KTSKEHVAKM 45 201 MKDLESLQTA EQSLMDLQER LEKAQEENRN VAVEKQNLER KMMDEINYAK 251 EEACRLRELR EGAECELSRR QYAEQELEQV RMALKKAEKE FELRSSWSVP 301 DALQKWLQLT HEVEVQYYNI KRQNAENQLA IAKDEAEKIK KKRSTVFGTL 351 HVAHSSSLDE VDHKILEAKK ALSELTTCLR ERLFRWQQIE KICGFQIAHN 50 401 SGLPSLTSSL YSDHSWVWP RVSIPPYPIA GGVDDLDEDT PPIVSQFPGT 451 MAKPPGSLAR SSSLCRSRRS IVPSSPQPQR AQLAPHAPHP SHPRHPHHPQ 501 HTPHSLPSPD PDILSVSSCP ALYRNEEEEE AIYFSAEK@W EVPDTASECD 551 SLNSSIGRKQ SPPLSLEIYQ TLSPRKISRD EVSLEDSSRG DSPVTVDVSW LOD GSPDCVGLTE TKSMIFSPAS KVYNGILEKS CSMNQLSSGI PVPKPRHTSC 55 P27 SZYBUDZKAN GEADZAKIZ ZIBHDFCHNG EKZKKDZKIK ZFŁKKKZK

BLASTP hits

No BLASTP hits available 5 Alert BLASTP hits for DKFZphamy2_24b4, frame 3 No Alert BLASTP hits found Pedant information for DKFZphamy2_24b4, frame 3 10 Report for DKFZphamy2_2464.3 ELENGTHD 769 EMWI 86673.49 [[q] 6.69 TREMBL:HS5242610_1 gene: "STIM1"; product: "GOK"; EHOMOLI Homo sapiens GOK (STIML) mRNA, complete cds. le-154 20 EBLOCKSI BLODBALC Dihydroxy-acid and L-phosphogluconate dehydratases proteins **EBFOCK2** PR00021D **EBFOCKZ** PR01053F EBF0CKZ] BLOO726B AP endonucleases family 1 proteins 25 EPROSITE RGD ľ BE BUTTEL JANGEZ [KW] EKWI TRANSMEMBRANE 1 [KW] LOW_COMPLEXITY 15.86 % [KW] COILED_COIL 8.45 % 30 SEQ RLHPASTPGPAWGWLLRRRRWAALLVLGLLVPGAADGCELVPRHLRGRRATGSAATAASS SEG PRD 35 COILS MEM PAAAAGDSPALMTDPCMSLSPPCFTEEDRFSLEALQTIHKQMDDDKDGGIEVEESDEFIR SEQ 40 SEG xxxxxxxx....... PRD ccccccccccccchhhhhhhhhhhhhhhhhhccccceeeecchhhhh COILZ MEM 45 EDMKYKDATNKHSHLHREDKHITIEDLWKRWKTSEVHNWTLEDTLQWLIEFVELPQYEKN SEQ SEG PRD COILZ 50 FRDNNVKGTTLPRIAVHEPSFMISQLKISDRSHRQKLQLKALDVVLFGPLTRPPHNWMKD SEQ SEG 55 PRD

COILS

MEM

5	SEG PRD COIL	hpheeeeeccccceeeecccccppppppppppppppppp
	MEM	MMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMM
10	SEQ SEG PRD COIL	
15	MEM	cccccccccccccccccccccccccccccccccccccc
	SEQ SEG PRD COIL:	
20	MEM	••••••••••••••••••••••••
25	SEQ SEG PRD COIL:	LHVAHSSSLDEVDHKILEAKKALSELTTCLRERLFRWQQIEKICGFQIAHNSGLPSLTSS eeeeeccccchhhhhhhhhhhhhhhhhhhhhhhhhh
	MEM	***************************************
30	SEQ SEG PRD COILS	LYSDHSWVWBRVSIPAYBAGVDDLDEDTDTDVBQFPGTMAKPPGSLARSSSLCRSRR CCCCCCCCCCCCCCCCCCCCCCCCCCCCC
35	MEM	***************************************
40	SEQ SEG PRD COILS	SIVPSSPQPQRAQLAPHAPHPHHPHHPHHPHDHTHSLPSPDPDILSVSSCPALYRNEEEE xxxxxxxxxxxxxxxxxxxxxxxxxxxxx
	MEM	
45	SEQ SEG PRD COILS	
50	MEM	
	SEQ SEG PRD	GDSPVTVDVSWGSPDCVGLTETKSMIFSPASKVYNGILEKSCSMNQLSSGIPVPKPRHTS
55	COILZ	CCCCeeeeeccccccccccccccccccccccccccccc
	MEM	
	SEQ	CZSAGNDZKPVQEAPZVARIZZIPHDLCHNGEKZKKPZKIKZLFKKKZK

DKFZphamy2_24c8

5 group: transmembrane protein

DKFZphamy2_24cA encodes a novel 454 amino acid protein without similarity to known proteins.

- The novel protein contains 1 transmembrane region.
 No informative BLAST results; No predictive prosite, pfam or SCOP motife.
- The new protein can find application in studying the expression profile of amygdala-specific genes and as a new marker for amygdala cells.
- 20 putative protein

EST of GEN-426HO7 is 141 Bp longer at 5'-end perhaps complete cds. Pedant: TRANSMEMBRANE 1

25 Sequenced by GBF

Locus: /map="609.7 cR from top of Chr3 linkage group"

30 Insert length: 3200 bp
Poly A stretch at pos. 3177, polyadenylation signal at pos. 3156

1 CCTGTCCACA GGGCCCGCTC CAGCAGCCAT GGCAACCACA TCCTCCAAGC 35 51 CAGAGGGCCG CCCTCGAGGG CAGGCTGCCC CCACCATCCT GCTGACAAAG JOI CCACCGGGGG CCACCAGCCG CCCCACCACA GCGCCCCCCC GCACTACCAC 151 ACGCAGGCCC CCCAGGCCCC CAGGCTCTTC CCGAAAAGGG GCTGGTAATT 201 CATCACGCCC TGTCCCGCCT GCACCTGGTG GCCACTCCAG GAGTAAAGAA 251 GGACAGCGAG GACGAAATCC AAGCTCCACA CCTCTGGGGC AGAAGCGGCC 40 301 CCTGGGGAAA ATCTTTCAGA TCTACAAGGG CAACTTCACA GGGTCTGTGG 351 AACCGGAGCC CTCTACCCTC ACCCCCAGGA CCCCACTCTG GGGCTACTCC 4D1 TCTTCACCAC AGCCCCAGAC AGTGGCTGCG ACCACAGTGC CCAGCAATAC 451 CTCATGGGCA CCCACCACCA CCTCCCTGGG GCCTGCAAAG GACAAGCCAG 5D1 GCCTTCGCAG AGCAGCCCAG GGGGGTGGTT CTACCTTCAC CAGCCAAGGA 45 551 GGGACACCAG ATGCCACAGC AGCCTCAGGT GCCCCTGTCA GTCCACAAGC LOD TGCCCCAGTG CCTTCTCAGC GCCCCCACCA CGGTGACCCA CAGGATGGCC **651 CCAGCCATAG TGACTCTTGG CTTACTGTTA CCCCTGGCAC CAGCAGACCT** 701 CTGTCTACCA GCTCTGGGGT CTTCACGGCT GCCACGGGGC CCACCCCAGC 751 TGCCTTCGAT ACCAGTGTCT CAGCCCCTTC CCAGGGGATT CCTCAGGGAG 50 BD1 CATCCACAAC CCCACAAGCT CCAACCCATC CCTCCAGGGT CTCAGAAAGC B51 ACTATTTCTG GAGCCAAGGA GGAGACTGTG GCCACCCTCA CCATGACCGA 901 CCGGGTGCCC AGTCCTCTCT CCACAGTGGT ATCCACAGCC ACAGGCAATT 951 TCCTCAACCG CCTGGTCCCC GCCGGGACCT GGAAGCCTGG GACAGCAGGG DOD AACATCTCCC ATGTGGCCGA GGGGGACAAA CCGCAGCACA GAGCCACCAT
DD51 CTGCCTGAGC AAGATGGATA TCGCCTGGGT GATCCTGGCC ATCAGCGTGC
LLD1 CCATCTCCTC CTGCTCTGTC CTGCTGACGG TGTGCTGCAT GAAGAGGAAG
LL51 AAGAAGACCG CCAACCCGGA GAACAACCTG AGCTACTGGA ACAACACCAT 55 1201 CACCATGGAC TACTTCAACA GGCATGCTGT GGAGCTGCCC AGGGAGATCC

WO 01/98454

PCT/IB01/02050

1251 AGTCCCTTGA AACCTCTGAG GACCAGCTCT CAGAGCCCCG CTCCCCAGCC 1301 AATGGCGACT ATAGAGACAC TGGGATGGTC CTTGTTAACC CCTTCTGTCA 1351 AGAAACACTG TTTGTGGGAA ACGATCAAGT ATCTGAGATC TAACTACAGC
1401 AGGCATCACT TTGCCATTCC GTATTTTTCG TCTCTAAATT ATAAATATAC 1451 AAATATATA ATTATAAATA TAACCTTTGT GTAACCCTGA CTTAATGAGA 5 1501 AACATTTCA GCTTTTTTC CTATGAATTG TCAACATCTT TTTTACAAGT 1551 GTGGTTTAAA AAAAAAAAA CTTTACAGAA TGATCTGTGG CTTTATAAAA
1601 TAAAGGTATT TCTAAGCAAA GCAGTTGCAT TGATTGCTTC TCTTAATAAC
1651 TATTCTTGAG CACCTGGGGA TCCCAGGAAC CCTGGTCAGG TGAGGTAAGA
1701 GACTGACCTC CTGTAGAAGC TGAATGTTAC AGTGGTCAAG CGCACGATTC
1751 TTTGAGTGAT TCTTAAAGCT CTGGTTCCTC TTGATTTGGT GTGACCCCAT
1801 TTCCTCCCTT CTCATACGCA CACCTGTAAA GGGAACTGGA CCGCCTCAGG
1851 GGAAGACGGC AGACTCATGC ACAGAGAAGG AAAAGGGAAC ATCTCATCAC
1951 GTTTAATTCC ATCCAAGTTG TGGATGGCAG GCAGGAGCAT GGAGCCCTCA
1951 GTTTAATTCC ATCCAAGTTG TGGATGGCAG GCAGGAGCAT GGAGCCCTCA
2001 GGAATCCATG GAGGACATCA AGGCCATCCCA AGGCCCATATT CCCCTAACAT 10 2001 GGAATCCATG GAGGACATCA AGGCATCCCA AGGCCATATT CCCCTAACAT 2051 TACTTCCACT GCTAACAACA GGACTGCCTT TCCCTGGTGG GAAAATGCTC
2101 CCTTTATGCC CATTCCTGTA TCCCCTCCAA CACCCACATC TGCATTAAAC
2151 ACCCGTGCCT TTCTCTTGGA GAGGGTTTAG ATGCAGATCC CGGCCCTGGA 22D1 GCTTTAAAAT GCTTGCCCTT CCTTCTTCAA GGATCAAATG TTTATTGGGG 20 2251 TTCAGCTTTG TTTTCTCAAA AGGCCATGGT ATCGTGCCCC TGAGGAACAT 2301 GTTATCTAA GAAGCTTTGA GGTAGTAGAG CGATAATTT TGAAACCTTC 2301 GTTTATCTAA GAAGCTTTGA GGTAGTAGAG CGATAATTTT TGAAACCTTC
2351 CTCCTGCAAT CTTTAAAAAA GAAAAAAAAG ATTGCCCAAA CAAATCATTT
2401 GGGAGAAGAC ATCATTATAC TCCTACTTGG CACTGCAAAC CTGCTCGCAG
2451 CACCAGCCGG TGGACTTGCC ATCCAGCTCT CAGCTTCCAC TGCTCCCCTT
2501 GTTCCCGGCC GGCTGGCTGC CTCCCCGTGC TGTGTCCAGC ACGGCCAACA
2551 ACGTCAGACC CTCAGAGACG CCCAAGGGGC TTCCAGAGGT GGCCGCTTCT
2601 CTATTTTTC CTGATTGTGG CTGAGAGAGA TGATTACTGC TTTGACACTT
2651 CCTTTCTCTA AAAGAAAAAA AGTTTGATAG TATATTTTGA ATATAGATGC
2701 TCTTATAGTC AGATTGGGAA TTGAACTTGA ATATTGGGTC ATATGTTTGT
2751 GTTGTTGCTG TAGTCTATCA TGACTTTTTT CTTTCTGCAT TTTCCTTAAA
2801 AAAAAAAAAA AGATGGCCTT CAAAAGTGTG TTCTCAATGT TGTATGAACC
2851 TCCTTCACAT GAGTTCGGTT GTTGTCTCTC TTCAAAGACC CTTCAACCCA
2901 CAAAGAAGCA ACTAAATGTT TCTCTAAGTT TAATTTTCTA GCGTGTTGTT
2951 GTCTTACCTT TTTAACCTTA CCATAATATT TCTGTTAACT GTTACATTTA
3001 ATATACCAAT GTGTGTAAGT ATACAGAGAA AAATCTGTTT GTAAAGGTAAA
3051 ATTTATATAT AATATATGTA ATCAAAAGATA CATATGTTAT ATATACATAT 30 3051 ATTTATAT AATATATGTA ATCAAAGATA CATATGTTAT ATATACATAT 3101 GTGGATGTAT GACTTATTTT TCCTTATCCA CAGATTTCAG CTACCATGTA 3151 TATAAATA AACTTATTTT ATTAGCCAGA GAAAAAAAA AAAAAAAAA 40

BLAST Results

45 No BLAST result

Medline entries

50 No Medline entry

55

Peptide information for frame 2

ORF from 29 bp to 1390 bp; peptide length: 454

Category: putative protein Classification: Transmembrane proteins unclassified 1 MATTSSKPEG RPRGQAAPTI LLTKPPGATS RPTTAPPRTT TRRPPRPPGS 5 51 SRKGAGNSSR PVPPAPGGHS RSKEGQRGRN PSSTPLGQKR PLGKIFQIYK JOI GNFTGSVEPE PSTLTPRTPL WGYSSSPQPQ TVAATTVPSN TSWAPTTTSL 151 GPAKDKPGLR RAAQGGGSTF TSQGGTPDAT AASGAPVSPQ AAPVPSQRPH 201 HGDPQDGPSH SDSWLTVTPG TSRPLSTSSG VFTAATGPTP AAFVFSWRPH 251 SQGIPQGAST TPQAPTHPSR VSESTISGAK EETVATLTMT DRVPSPLSTV 301 VSTATGNFLN RLVPAGTWKP GTAGNISHVA EGDKPQHRAT ICLSKMDIAW 351 VILAISVPIS SCSVLLTVCC MKKKKKTANP ENNLSYWNNT ITMOYFNRHA 10 401 VELPREIQSL ETSEDQLSEP RSPANGDYRD TGMVLVNPFC QETLFVGNDQ 451 VSEI 15 BLASTP hits. No BLASTP hits available 20 Alert BLASTP hits for DKFZphamy2_24c8, frame 2 No Alert BLASTP hits found Pedant information for DKFZphamy2_24c8, frame 2 25 Report for DKFZphamy2_24c8-2 30 ELENGTHD 463 48277-84 [[q] 9-80 EFUNCATI 98 classification not yet clear-cut ES. cerevisiae = 35 YJR151c3 2e-04 **EBFOCK23** PR00912F EBLOCKS3 BPO3696F EKW3 TRANSMEMBRANE EKWI 15.55 % LOW_COMPLEXITY 40 SEQ LSTGPAPAAMATTSSKPEGRPRGQAAPTILLTKPPGATSRPTTAPPRTTTRRPPRPPGSS SEG PRD 45 MEM RKGAGNSSRPVPPAPGGHSRSKEGQRGRNPSSTPLGQKRPLGKIFQIYKGNFTGSVEPEP SEQ SEG PRD 50 MEM STLTPRTPLWGYSSSPQPQTVAATTVPSNTSWAPTTTSLGPAKDKPGLRRAAQGGGSTFT SEQ SEG PRD 55 MEM

ZQGGTPDATAASGAPVSPQAAPVPSQRPHHGDPQDGPSHSDSWLTVTPGTSRPLSTSSGV

SEQ

SEG

	W	WO 01/98454	PCT/IB01/02050					
	PRD MEM							
5	SEQ SEG PRD MEM	G xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx	eeeccchhhhhhhccc					
10	SEQ SEG PRD MEM	G D cccccceeeeecccccccccccccccccceeecccc	ccceeecccchhhhhh					
15	SEQ SEG PRD MEM	hhhhccccccceeeehhhhhhccccccccccccccc	ccccccccchhhhhc					
20	SEQ SEG PRD MEM	G D cccccccccccccccceeeeeccccceeeeecccccc	•					
25	(No	o Prosite data available for DKFZphamy2_24c8.	2)					
	(No	o Pfam data available for DKFZphamy2_24c8.2)						

WO 01/98454



DKFZphamy2_24k15

5 group: amygdala derived

DKFZphamy2_24kl5 encodes a novel 279 amino acid protein with weak similarity to pecanex of Drosophila melanogaster-

- 10 Pecanex is a maternal-effect neurogenic gene involved in differentiation processes in the developing central nervous system. DKFZphamy2_24kl5.p3 seems to be expressed ubiquitiously.
- The new protein can find application in studying the expression profile of amygdala-specific genes and as a new marker for amygdala cells.

similarity to pecanex (Drosophila melanogaster)

20 probably complete cds.

Sequenced by GBF

25 Locus: unknown

Insert length: 1464 bp

Poly A stretch at pos. 1445, polyadenylation signal at pos. 1421

30 L AAGGAAAACA AGAGGACATG CCATATATTC CTCTCATGGA GTTCAGTTGT 51 TCACATTCTC ACTTAGTATG CTTACCCGCA GAGTGGAGGA CTAGCTGTAT JOB GCCCAGTTCC AAAATGAAGG AGATGAGCTC GTTATTTCCA GAAGACTGGT 151 ACCAATITGT TCTAAGGCAG TTGGAATGTT ATCATTCAGA AGAGAAGGCC 201 TCAAATGTAC TGGAAGAAAT TGCCAAGGAC AAAGTTTTAA AAGACTTTTA 35 251 TGTTCATACA GTAATGACTT GTTATTTTAG TTTATTTGGA ATAGACAATA 3D1 TGGCTCCTAG TCCTGGTCAT ATATTGAGAG TTTACGGTGG TGTTTTGCCT 351 TGGTCTGTTG CTTTGGACTG GCTCACAGAA AAGCCAGAAC TGTTTCAACT HOL AGCACTGAAA GCATTCAGGT ATACTCTGAA ACTAATGATT GATAAAGCAA 40 451 GTTTAGGTCC AATAGAAGAC TTTAGAGAAC TGATTAAGTA CCTTGAAGAA 501 TATGAACGTG ACTGGTACAT TGGTTTGGTA TCTGATGAAA AGTGGAAGGA 551 AGCAATTTTA CAAGAAAAGC CATACTTGTT TTCTCTGGGG TATGATTCTA LOI ATATGGGAAT TTACACTGGG AGAGTGCTTA GCCTTCAAGA ATTATTGATC **L51 CAAGTGGGAA AGTTAAATCC TGAAGCTGTT AGAGGTCAGT GGGCCAATCT** 45 701 TTCATGGGAA TTACTTTATG CCACAAACGA TGATGAAGAA CGTTATAGTA 751 TACAAGCTCA TCCACTACTT TTAAGAAATC TTACGGTACA AGCAGCAGAA BOL CCTCCCTGG GATATCCGAT TTATTCTTCA AAACCTCTCC ACATACATTT 851 GTATTAGAGC TCATTTTGAC TGTAATGTCA TCAAATGCAA TGTTTTTATT 901 TTTTCATCCT AAAAAAGTAA CTGTGATTCT TGTAACTTGA GGACTTCTCC 50 95% ACACCCCCAT TCAGATGCCT GAGAACAGCT AAGCTCCGTA AAGTTGGTTC DOD TCTTAGCCAT CTTAATGGTT CTAAAAAACA GCAAAAACAT CTTTATGTCT 1051 AAGATAAAAG AACTATTTGG CCAATATTTG TGCCCTCTGG ACTTTAGTAG 1303 GCTTTGGTAA ATGTGAGAAA ACTTTTGTAG AATTATCATA TAATGAATTT 1151 TGTAATGCTT TCTTAAATGT GTTATAGGTG AATTGCCATA CAAAGTTAAC 1201 AGCTATGTAA TTTTTACATA CTTAAGAGAT AAACATATCA GTGTTCTAAG 55 1251 TAGTGATAAT GGATCCTGTT GAAGGTTAAC ATAATGTGTA TATATTTGTT 1301 TGAAATAA TTTATAGTAT TTTCAAATGT GCTGATTTAT TTTGACATCT 1351 AATATCTGAA TGTTTTTGTA TCAAGTAGTT TGTTTTCATA GACTTCAATT

5 BLAST Results

Entry ACOO7939 from database EMBLNEW:
Homo sapiens clone 422_H_5, WORKING DRAFT SEQUENCE, 5 unordered
pieces.
Score = 4116, P = 0.0e+00, identities = 840/858
3 exons

.

15

Medline entries

No Medline entry

20

Peptide information for frame 3

25

ORF from 18 bp to 854 bp; peptide length: 279 Category: similarity to known protein Classification: unclassified

30 J MPYIPLMEFS CSHSHLVCLP AEWRTSCMPS SKMKEMSSLF PEDWYQFVLR
51 QLECYHSEEK ASNVLEEIAK DKVLKDFYVH TVMTCYFSLF GIDNMAPSPG
101 HILRVYGGVL PWSVALDWLT EKPELFQLAL KAFRYTLKLM IDKASLGPIE
151 DFRELIKYLE EYERDWYIGL VSDEKWKEAI LQEKPYLFSL GYDSNMGIYT
201 GRVLSLQELL IQVGKLNPEA VRGQWANLSW ELLYATNDDE ERYSIQAHPL
35 251 LLRNLTVQAA EPPLGYPIYS SKPLHIHLY

BLASTP hits

40

55

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_24kl5, frame 3

45 No Alert BLASTP hits found

Pedant information for DKFZphamy2_24kl5, frame 3

Report for DKFZphamy2_24kl5.3

 ELENGTHI
 284

 EMUI
 33066.31

 Epii
 5.17

TREMBL: AFO67608_11 gene: "BO511.12";
Caenorhabditis elegans cosmid BO511.2e-13

EKWI Alpha_Beta

5	SEQ PRD	GKQEDMPYIPLMEFSCSHSHLVCLPAEWRTSCMPSSKMKEMSSLFPEDWYQFVLRQLECY CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
	SEQ PRD	hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
10	SEQ PRD	LDWLTEKPELFQLALKAFRYTLKLMIDKASLGPIEDFRELIKYLEEYERDWYIGLVSDEK cchhhhhhchhhhhhhhhhhhhhhhhhhhhhhhhhhh
	SEQ PRD	WKEAIL@EKPYLFSLGYDSNMGIYTGRVLSL@ELLI@VGKLNPEAVRG@WANLSWELLYA
15	SEQ PRD	TNDDEERYSIQAHPLLLRNLTVQAAEPPLGYPIYSSKPLHIHLY CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
20	(No	Prosite data available for DKFZphamy2_24k15.3)
	(No	Pfam data available for DKFZphamy2_24k35.3)

DKFZphamy2_2al3

5 group: amygdala derived

DKFZphamy2_2al3 encodes a novel 440 amino acid protein without similarity to known proteins.

No informative BLAST results; No predictive prosite, pfam or SCOP motife.

The new protein can find application in studying the expression profile of amygdala-specific genes.

15

putative protein

perhaps complete cds.

20

Sequenced by MediGenomix

Locus: /map="16p13.3"

25 Insert length: 2584 bp

Poly A stretch at pos. 2562, polyadenylation signal at pos. 2545

1 GTTCCTGAGG ACGTGCTACG GGGGCAGCTT CCTGGTACAC GAGTCGTTCC

```
51 TCTACAAGCG GGAGAAGGCT GTCGGGGACA AGGTGTATTG GACCTGCCGG
101 GACCACGCGC TGCACGGCTG CCGGAGCCGG GCCATCACCC AGGGACAGCG
30
           151 GGTGACTGTG ATGCGTGGGC ACTGCCACCA GCCCGATATG GAGGGCCTGG
           201 AAGCCCGGCG GCAGCAGGAG AAGGCCGTGG AGACGCTGCA GGCTGGGCAG
           251 GACGGCCCTG GGAGCCAAGT GGACACGCTG CTCCGAGGCG TGGATAGTTT
           ADD GCTCTACCGC AGGGGTCCGG GTCCCCTGAC TCTCACCAGG CCTCGGCCCA
35
           351 GAAAGCGAGC AAAGGTCGAA GACCAGGACC TGCCAACCCA GCCCGAGGCC
           4D1 CCAGACGAGC ACCAGGACAT GGACGCAGAC CCGGGAGGCC CTGAGTTCCT
           451 GAAGACGCCC CTGGGGGGCA GCTTCCTGGT GTACGAGTCC TTCCTCTACC
           501 GGCGGGAGAA GGCGGCTGGG GAGAAGGTGT ATTGGACCTG CCGGGACCAG
40
           551 GCCCGCATGG GCTGCCGCAG CCGCGCCATC ACCCAGGGCC GACGGGTGAC
           LOD TGTCATGCGT GGTCACTGCC ACCCGCCCGA CCTGGGAGGC CTGGAGGCCC
           LSD TGAGGCAGCG GGAGAAACGC CCCAACACGG CGCAGCGGGG GAGCCCAGGC 7DD GCTGGCCTCT CTTTCCAGTG GCTCTTCCGG ATCCTGCAGC TTTTGGGTCA
          751 TGCTCCTGTG CTGCTGTGCC CCTCAGGGTC CTCCTGCCTC CCGAGCCTCC
BD1 CTGCTCCACA TGGCCCCTGC CCAGCCCTCT CCATCCCTCT TGAAGGAGGC
B51 CCCGAGTTCC TGAAGACGCC CCTGGGGGGC AGCTTCCTGG TGTACGAGTC
45
           POL CTTCCTCTAC CGGCGGGAGA AGGCGGCCGG GGAGAAGGTG TATTGGACCT
           951 GCCGGGACCA GGCCCGCATG GGCTGCCGCA GCCGCGCCAT CACCCAGGGC
        TOTAL GCCGGGACCA GGCCCGCATG GGCTGCCGCA GCCGCGCCAT CACCCAGGGC

LODI CGGCGGGTCA TGGTCATGCG CAGGCACTGC CACCCACCGG ACCTGGGCGG

LODI CCTGGAGGCC CTGCGGCAGC GGGAGCACTT CCCCAACCTG GCGCAGTGGG

LIDI ACAGCCCAGA TCCTCTCCGG CCCCTGGAGT TCCTGAGGAC TTCCCTGGGG

LIDI TGGGGAGAAG GTGTACTGGA TGTGCCGGGA CCAGGCTCGG CTGGGCTGCC

LODI TGGGGAGAAG GTGTACTGGA TGTGCCGGGA CCAGGCTCGG CTGGGCTGCC

LODI TGCCATCAGC CATAACCCAG GGCCACCGCA TCATGGTCAT GCGCAGCCAC

LODI TGCCATCAGC CTGACCTGGC AGGCCTGGAG GCCTTGAGGC AACGGGAGCG

LODI TGCCATCAGC CTGACCTGGC AGGAGGACCC AGAAAAGATT CAAGTTCAGC

LODI TGCCATCAAC CTGACCTGGC AGGAGGACCC AGAAAAGATT CAAGTTCAGC

LODI TGTGCTTCAA GACGTGTTCT CCTGAAAGCC AGCAGATTTA TGGGGACATTC
50
55
         1401 TGTGCTTCAA GACGTGTTCT CCTGAAAGCC AGCAGATTTA TGGGGACATC
        1451 AAAGACGTCA GACTGGATGG CGAGTCCCAG TGAGGCGATG TGGGCAGAGG
```

```
1501 AGCTCCGAGC CGCCCACCCA AGGTGGCTTC ACATCCACAC AGGCACTTCC
      1551 CATCCACCTA GGTTTGGCTT AGCAGAAACT TCTTTTCATT CTTCCAAAGC
      1601 ATCGATGGTC TTCGCGTCTC CTCAGGAGGT CTCCCAGGAG GAATTCTTGG
      1651 ATGGTGTCCT CATGTCGGCG GAGAACAGTG CTCAGAGCTG GCGCTTGCAG
      1701 ACGCAGCTGT CGTGGGGCAG GGCGGTGGCG CCTTCCTGAC CTTTGGAAGA
5
      1751 CATGACAAAG CTGCCTGGAC ACGGACGCCC CTGCTGTACG GCCACAGCAC
      LBOL CCCTGGGTTT GCAGAGCACG CAGCCTTCCT AGGGCTTTCC ACCTGGCGAG
      1851 GCCCCGCTCT GCTCAGCACG GTGCAAAGTG AATGCTGCTG TCTTGGAGCC
      1901 TGGGCACGTT TGGGGAAGTT CCTGCTTCAA ACTGAGCTGC CCCGCATAGG
      1951 CCAGGTCAAC CCACACCAAT CTTTTCTGGA CAGGTGCTGG GTAGGCCTTC
10
      2001 CTGGTCTCTG GCCGCCTGCT GCCAGGGTGT GGCCATCCCC AGCAACCGGA
      2051 GCCGGCCAAA CCAGAGGCCT CGCTCCGCAC TCCACACTTT CCTTTCTGTG
      PLO CTCCTTCCAA GTTAAATTAA ACCCCCTCTC CACGATTCCC ACGGCAGGCG
      2151 TCATTCCCGA GATGGGAGCC AGTCCAGGGG TCAGCAGGAG CCAGCGCTGG
      2201 GCACACGTGC CCTGGCTGAG GCCAGCGGCA TCCTGGGTGG CCCAGGTCCA
15
      2251 TCCTGGCAG CAAAGGCGTG TCCCCTTCTG TCAGACAGCT TCACAGAGTG
2301 TGGCTTCACC AGTCAGAGGG AGCAGTCCGG AGAGGCAAGA TGACCCCACC
2351 GGGACTGCAG AGCCTCCTCC TTACTAACAA GGACCTGTCC GCAGCCGCGA
2401 GGTCCTTCAC TCCCACCCTG TAATTGTGGG GGGAGTGCCA GCAACAGGCC
2451 TGTCCCCTGG CAAGTTGGCC ACGGAACCCA CCATGCACTG CAAGGCTGTG
2501 ACAGCCTTGGG CACCCCTGCT TCTCCTCTGC TTACTAACGGTT CCCCCCAATAA
20
      2551 ATCCTATTTT CCATCAAAAA AAAAAAAAAA AAAA
```

25 BLAST Results

No BLAST result

30

Medline entries

No Medline entry

35

Peptide information for frame 2

40

ORF from 161 bp to 1480 bp; peptide length: 440 Category: putative protein Classification: no clue

1 MRGHCHQPDM EGLEARRQQE KAVETLQAGQ DGPGSQVDTL LRGVDSLLYR
51 RGPGPLTLTR PRPRKRAKVE DQELPTQPEA PDEHQDMDAD PGGPEFLKTP
101 LGGSFLVYES FLYRREKAAG EKVYWTCRDQ ARMGCRSRAI TQGRRVTVMR
151 GHCHPPDLGG LEALRQREKR PNTAQRGSPG AGLSFQWLFR ILQLLGHAPV
201 LLCPSGSSCL PSLPAPHGPC PALSIPLEGG PEFLKTPLGG SFLVYESFLY
50 251 RREKAAGEKV YWTCRDQARM GCRSRAITQG RRVMVMRRHC HPPDLGGLEA
301 LRQREHFPNL AQWDSPDPLR PLEFLRTSLG GRFLVHESFL YRKEKAAGEK
351 VYWMCRDQAR LGCRSRAITQ GHRIMVMRSH CHQPDLAGLE ALRQRERLPT
401 TAQQEDPEKI QVQLCFKTCS PESQQIYGDI KDVRLDGESQ

55

BLASTP hits

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_2al3, frame 2

No Alert BLASTP hits found

493

Pedant information for DKFZphamy2_2al3, frame 2

10 Report for DKFZphamy2_2al3.2

ELENGTHI EMMI 55840.13 9.33 15 [pI] Alpha_Beta [KW] [KW] LOW_COMPLEXITY 6.29 % 20 SEQ FLRTCYGGSFLVHESFLYKREKAVGDKVYWTCRDHALHGCRSRAITQGQRVTVMRGHCHQ SEG PRD SEQ PDMEGLEARRQQEKAVETLQAGQDGPGSQVDTLLRGVDSLLYRRGPGPLTLTRPRPRKRA 25 SEG PRD cccchhhhhhhhhhhhhhhhccccccccccccccceeeecccceehhh KVEDQELPTQPEAPDEHQDMDADPGGPEFLKTPLGGSFLVYESFLYRREKAAGEKVYWTC SEQ SEG 30 PRD SEQ RDQARMGCRSRAITQGRRVTVMRGHCHPPDLGGLEALRQREKRPNTAQRGSPGAGLSFQW SEG PRD 35 SEQ LFRILQLLGHAPVLLCPSGSSCLPSLPAPHGPCPALSIPLEGGPEFLKTPLGGSFLVYES SEG PRD 40 SEQ FLYRREKAAGEKVYWTCRDQARMGCRSRAITQGRRVMVMRRHCHPPDLGGLEALRQREHF SEG PRD SEQ PNLAQWDSPDPLRPLEFLRTSLGGRFLVHESFLYRKEKAAGEKVYWMCRDQARLGCRSRA 45 SEG PRD SEQ ITQGHRIMVMRSHCHQPDLAGLEALRQRERLPTTAQQEDPEKIQVQLCFKTCSPESQQIY SEG 50 SEQ GDIKDVRLDGESQ ZEG PRD ccccccccc 55

(No Prosite data available for DKFZphamy2_2al3.2)

(No Pfam data available for DKFZphamy2_2al3.2)

DKFZphamy2_2bl9

5 group: differentiation/development

DKFZphamy2_2bl9 encodes a novel 789 amino acid protein which originates from TXBPl51 mRNA by alternative splicing.

10

It is ubiquitously expressed. The mRNA is also subject to alternative polyadenylation. Overexpression of TXBPl5L in NIH3T3 cells causes inhibition of apoptosis induced by tumour necrosis factor (TNF). It binds to A2D, which is

15 A20, which is also an inhibitor of cell death by a yet unknown mechanism.

The new protein can find application in modifying/blocking apoptosic pathways and therefore serve as a tool in diagnosis of cancer predisposition and as a tool in cell culture.

TXBPL51 differentially spliced

25 differential splicing differential polyadenylation

Sequenced by MediGenomix

30 Locus: /map="7pl5"

Insert length: 3028 bp Poly A stretch at pos. 2885, polyadenylation signal at pos. 2868

35 1 GAAGAGGTTC GGCGGCTGAT GGCGGATCAG GATCGGAAGC CTGCGTAACT 51 TTCTCCCTTG ATCCGGGAGT CTTTCCACTG GATTCACAAT GACATCCTTT DOD CAAGAAGTCC CATTGCAGAC TTCCAACTTT GCCCATGTCA TCTTTCAAAA 151 TGTGGCCAAG AGTTACCTTC CTAATGCACA CCTGGAATGT CATTACACCT 201 TAACTCCATA TATTCATCCA CATCCAAAAG ATTGGGTTGG TATATTCAAG 251 GTTGGATGGA GTACTGCTCG TGATTATTAC ACGTTTTTAT GGTCCCCTAT 301 GCCTGAACAT TATGTGGAAG GATCAACAGT CAATTGTGTA CTAGCATTCC 40 351 AAGGATATTA CCTTCCAAAT GATGATGGAG AATTTTATCA GTTCTGTTAC 4D1 GTTACCCATA AGGGTGAAAT TCGTGGAGCA AGTACACCTT TCCAGTTTCG 451 AGCTTCTTCT CCAGTTGAAG AGCTGCTTAC TATGGAAGAT GAAGGAAATT
501 CTGACATGTT AGTGGTGACC ACAAAAGCAG GCCTTCTTGA GTTGAAAATT
551 GAGAAAACCA TGAAAGAAAA AGAAGAACTG TTAAAGTTAA TTGCCGTTCT 45 LDL GGAAAAGAA ACAGCACAAC TTCGAGAACA AGTTGGGAGA ATGGAAAGAG 651 AACTTAACCA TGAGAAAGAA AGATGTGACC AACTGCAAGC AGAACAAAAG 701 GGTCTTACTG AAGTAACACA AAGCTTAAAA ATGGAAAATG AAGAGTTTAA
751 GAAGAGGTTC AGTGATGCTA CATCCAAAGC CCATCAGCTT GAGGAAGATA
801 TTGTGTCAGT AACACATAAA GCAATTGAAA AAGAAACCGA ATTAGACAGT 50 B51 TTAAAGGACA AACTCAAGAA GGCACAACAT GAAAGAGAAC AACTTGAATG 901 TCAGTTGAAG ACAGAGAAGG ATGAAAAGGA ACTTTATAAG GTACATTTGA 55 951 AGAATACAGA AATAGAAAAT ACCAAGCTTA TGTCAGAGGT CCAGACTTTA BOOL AAAAATTTAG ATGGGAACAA AGAAAGCGTG ATTACTCATT TCAAAGAAGA 1051 GATTGGCAGG CTGCAGTTAT GTTTGGCTGA AAAGGAAAAT CTGCAAAGAA 1101 CTTTCCTGCT TACAACCTCA AGTAAAGAAG ATACTTGTTT TTTAAAGGAG

1151 CAACTTCGTA AAGCAGAGGA ACAGGTTCAG GCAACTCGGC AAGAAGTTGT 1201 CTTTCTGGCT AAAGAACTCA GTGATGCTGT CAACGTACGA GACAGAACGA 1251 TGGCAGACCT GCATACTGCA CGCTTGGAAA ACGAGAAAGT GAAAAAGCAG ASTADAAAA ADTATSOTAA ACTTAAACTA AATGCTATGA AAAAAGATCA 1351 GGACAAGACT GATACACTGG AACACGAACT AAGAAGAGAA GTTGAAGATC 5 ኔዛዐኔ TGAAACTCCG TCTTCAGATG GCTGCAGACC ATTATAAAGA AAAATTTAAG 1451 GAATGCCAAA GGCTCCAAAA ACAAATAAAC AAACTTTCAG ATCAATCAGC 1501 TAATAATAAT AATGTCTTCA CAAAGAAAC GGGGAATCAG CAGAAAGTGA **3553 ATGATGCTTC AGTAAACACA GACCCAGCCA CTTCTGCCTC TACTGTAGAT** BEOD GTAAAGCCAT CACCTTCTGC AGCAGAGGCA GATTTTGACA TAGTAACAAA 10 1651 GGGGCAAGTC TGTGAAATGA CCAAAGAAAT TGCTGACAAA ACAGAAAAGT 1701 ATAATAAATG TAAACAACTC TTGCAGGATG AGAAAGCAAA ATGCAATAAA 1751 TATGCTGATG AACTTGCAAA AATGGAGCTG AAATGGAAAG AACAAGTGAA
1801 AATTGCTGAA AATGTAAAAC TTGAACTAGC TGAAGTACAG GACAATTATA LBOL AATTGCTGAA AATGTAAAAC TTGAACTAGC TGAAGTACAG GACAATTATA
LBSL AAGAACTTAA AAGGAGTCTA GAAAATCCAG CAGAAAGGAA AATGGAAGGT
LPOL CAGAATTCCC AGAGTCCTCA ATGTTTCAAA ACATGCTCAG AGCAAAATGG
LPSL TTATGTTCTC ACATTGTCAA ATGCACAACC AGTTCTGCAA TATGGTAATC
LPOL CTTATGCATC TCAGGAAACA AGAGATGGAG CAGATGGTGC TTTTTACCCA
LPSL GATGAAATAC AAAGGCCACC TGTCAGAAGT CCCTCTTGGG GACTGGAAGA
LPOL CAATGTTGTC TGCAGCCAGC CTGCTCGAAA CTTTAGTCGG CCTGATGGCT
LSSL TAGAGGACTC TGAGGATAGC AAAGAAGATG AGAATGTGCC TACTGCTCCT
LPOL GATCCTCCAA GTCAACATTT ACGTGGGCAT GGGACAGGCT TTTGCTTTGA
LPSL TTCCAGCTTT GATGTTCACA AGAAGTGTCC CCTCTGTGAG TTAATGTTTC
LPOL CTCCTAACTA TGATCAGAGC AAATTTGAAG AACATGTTGA AAGTCACTGG
LPSL AAGGTGTGCC CGATGTGCAG CGAGCAGTTC CCTCCTGACT ATGACCAGCA
LPOL GGTGTTTGAA AGGCATGTGC AGACCCATTT TGATCAGAAT GTTCTAAATT
LPSL TTGACTAGTT ACTTTTTATT ATGAGTTAAT ATAGTTTAGC AGTAAAAAAA
LSOL AAAAAAAAAA ACCACACCTA AAATAGACCA CTGAGGAGAC CATAGAGCGG 15 20 25 2501 AAAAAAAAA ACCACACCTA AAATAGACCA CTGAGGAGAC CATAGAGCGG 2551 ATGCTTTCAT GCACCCTTTA CTGCACTTTC TGACCAGGAG CTACTTTGAG 2601 TTTGGTGTTA CTAGGATCAG GGTCAGTCTT TGGCTTATCA ATAAATTTTA 30 2651 ATCTCTGTTA ATCTTACCTG CTTTAAAAAA AAGTTCTTGT GTGTTCGTAT 2701 CTTTATTTAT TCCCTAGTTT GCAGAACTGT CTGAATAAAG GATACAAGGA 2751 TTATTTCAAT GTTACTGCAC TGAAAAACGT GTATGTATTA GTGTGCTAGA 2803 TTATTTAGCA GAATATTCAC AAGTTTCTGT TGACCTTGTT GATTGAGCAT 2851 GACTACTAAA TATTATGTAA TAAAAAGCAT TTGTCATAAC AAAAAAAAA 35 AAAAAAA AAAAAAAA LODE

40

BLAST Results

No BLAST result

45

Medline entries

50 99361984: De Valck D, Jin DY, Heyninck K, Van de Craen M, Contreras R, Jeang KT, Beyaert R.; The zinc finger protein A2O interacts with 55

anti-apoptotic protein which is cleaved by specific caspases. Oncogene

1999 Jul 22:18(29):4182-90

5 Peptide information for frame 2 ORF from 89 bp to 2455 bp; peptide length: 789 Category: known protein Classification: Cell division 10 J MTSFQEVPLQ TSNFAHVIFQ NVAKSYLPNA HLECHYTLTP YIHPHPKDWV 51 GIFKVGWSTA RDYYTFLWSP MPEHYVEGST VNCVLAFQGY YLPNDDGEFY IDI @FCYVTHKGE IRGASTPF@F RASSPVEELL TMEDEGNSDM LVVTTKAGLL 151 ELKIEKTMKE KEELLKLIAV LEKETAQLRE QVGRMERELN HEKERCDQLQ 15 201 AEGKGLTEVT QSLKMENEEF KKRFSDATSK AHGLEEDIVS VTHKAIEKET 251 ELDSLKDKLK KARHERERLE CALKTEKDEK ELYKVHLKNT EIENTKLMSE 3D1 VQTLKNLDGN KESVITHFKE EIGRLQLCLA EKENLQRTFL LTTSSKEDTC 351 FLKEQLRKAE EQVQATRQEV VFLAKELSDA VNVRDRTMAD LHTARLENEK 403 VKKQLADAVA ELKLNAMKKD QDKTDTLEHE LRREVEDLKL RLQMAADHYK 20 451 EKFKECQRLQ KQINKLSDQS ANNNNVFTKK TGNQQKVNDA SVNTDPATSA 501 STVDVKPSPS AAEADFDIVT KGQVCEMTKE IADKTEKYNK CKQLLQDEKA 551 KCNKYADELA KMELKWKEQV KIAENVKLEL AEVQDNYKEL KRSLENPAER POT KWEGGNZGZP GCFKTCZEGN GYVLTLZNAG PVLGAGNPA ZGETRDGADG L51 AFYPDEIQRP PVRVPSWGLE DNVVCSQPAR NFSRPDGLED SEDSKEDENV 25 701 PTAPDPPSQH LRGHGTGFCF DSSFDVHKKC PLCELMFPPN YDQSKFEEHV 751 ESHWKVCPMC SEQFPPDYDQ QVFERHVQTH FDQNVLNFD 30 BLASTP hits No BLASTP hits available Alert BLASTP hits for DKFZphamy2_2bl9, frame 2 35 TREMBL:HS338211_1 product: "tax1-binding protein TXBP151"; Homo sapiens tax1-binding protein TXBP151 mRNA, complete cds., N = 2, Score = 2948, P = 0 40 >TREMBL:HS338211_1 product: "tax1-binding protein TX8P151"; Homo sapiens 45 taxl-binding protein TXBPL51 mRNA, complete cds. **Length** = 747 HSPs: Score = 2948 (442.3 bits), Expect = 0.0e+00, Sum P(2) = 0.0e+00 50 Identities = 575/603 (95%), Positives = 576/603 (95%) Querv: 55 MTSF@EVPL@TSNFAHVIF@NVAKSYLPNAHLECHYTLTPYIHPHPKDWVGIFKVGWSTA Sbjct: MTŠFQEVPLQTSNFAHVIFQNVAKSYLPNAHLECHYTLTPYIHPHPKDWVGIFKVGWSTA 60

	RDYYTFLWSPMPEHYVEGSTVNCVLAF@GYYLPNDDGEFY@FCYVTHKGEIRGASTPF@F	750
5	RDYYTFLWSPMPEHYVEGSTVNCVLAF@GYYLPNDDGEFY@FCYVTHKGEIRGASTPF@FSbjct: 61	
	RDYYTFLWSPMPEHYVEGSTVNCVLAFQGYYLPNDDGEFYQFCYVTHKGEIRGASTPFQF	750
10	Query: 121 RASSPVEELLTMEDEGNSDMLVVTTKAGXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX	180
	TAQLRE Sbjct: 121	
15	RASSPVEELLTMEDEGNSDMLVVTTKAGLLELKIEKTMKEKEELLKLIAVLEKETAQLRE	180
	QUERY: 181 QVGRMERELNHEKERCDQLQAEQKGLTEVTQSLKMENEEFKKRFSDATSKAHQLEEDIVS QVGRMERELNHEKERCDQLQAEQKGLTEVTQSLKMENEEFKKRFSDA	
20	+EEDIVS Sbjct: lal	
	QVGRMERELNHEKERCDQLQAEQKGLTEVTQSLKMENEEFKKRFSDATSKAHHVEEDIVS	240
25	VTHKAIEKETELDSLKDKLKKAQHEREQLECQLKTEKDEKELYKVHLKNTEIENTKLMSE	300
	VTHKAIEKETELDSLKDKLKKAQHEREQLECQLKTEKDEKELYKVHLKNTEIENTKLMSE Sbjct: 241	
	VTHKAIEKETELDSLKDKLKKAQHEREQLECQLKTEKDEKELYKVHLKNTEIENTKLMSE	300
30	QUERY: 301 VQTLKNLDGNKESVITHFKEEIGRLQLCLAEKENLQRTFLLTTSSKEDTCFLKEQLRKAE	360
	VQTLKNLDGNKESVITHFKEEIGRLQLCLAEKENLQRTFLLTTSSKEDTCFLKEQLRKAE Sbjct: 301	•
35	VQTLKNLDGNKESVITHFKEEIGRLQLCLAEKENLQRTFLLTTSSKEDTCFLKEQLRKAE	360
	Query: 361 EQVATRQEVVFLAKELSDAVNVRDRTMADLHTARLENEKVKKQLADAVAELKLNAMKKD	420
40	EQVQATRQEVVFLAKELSDAVNVADATHIDAMTRDAVAELKLNAMKKD Sbjct: 363	
	EQVQATRQEVVFLAKELSDAVNVRDRTMADLHTARLENEKVKKQLADAVAELKLNAMKKD	420
45	QUETY: 421 QDKTDTLEHELRREVEDLKLRLQMAADHYKEKFKECQRLQKQINKLSDQSANNNNVFTKK	480
	QDKTDTLEHELRREVEDLKLRLQMAADHYKEKFKECQRLQKQINKLSDQSANNNNVFTKK Sbjct: 421	
50	QDKTDTLEHELRREVEDLKLRLQMAADHYKEKFKECQRLQKQINKLSDQSANNNNVFTKK	480
	Query: 481 TGNQQKVNDASVNTDPATSASTVDVKPSPSAAEADFDIVTKGQVCEMTKEIADKTEKYNK	540
55	TGNQQKVNDASVNTDPATSASTVDVKPSPSAAEADFDIVTKGQVCEMTKEIADKTEKYNK Sbjct: 481	
-	TGNQQKVNDASVNTDPATSASTVDVKPSPSAAEADFDIVTKGQVCEMTKEIADKTEKYNK	540

Query: 541

CKQLLQDEKAKCNKYADELAKMELKWKEQVKIAENVKLELAEVQDNYKELKRSLENPAER 600

CKQLLQDEKAKCNKYADELAKMELKWKEQVKIAENVKLELAEVQDNYKELKRSLENPAER

5 Sbjct: 541

30

35

CKALLADEKAKCNKYADELAKMELKWKEAVKIAENVKLELAEVADNYKELKRSLENPAER POO

Query: 601 KME 603

KME

10 Sbjct: 601 KME 603

Score = 831 (124.7 bits), Expect = 0.0e+00, Sum P(2) = 0.0e+00 Identities = 147/153 (96%), Positives = 149/153 (97%)

15 Query: 637
NPASQETRDGADGAFYPDEIQRPVRVPSWGLEDNVVCSQPARNFSRPDGLEDSKE 696
NP A ++
DGADGAFYPDEIQRPVRVPSWGLEDVVVCSQPARNFSRPDGLEDSKE
Sbjct: 596 NP-

20 AERKMEDGADGAFYPDEIQRPPVRVPSWGLEDNVVCSQPARNFSRPDGLEDSEDSKE 654

Query: 697
DENVPTAPDPPSQHLRGHGTGFCFDSSFDVHKKCPLCELMFPPNYDQSKFEEHVESHWKV 756

25 DENVPTAPDPPSQHLRGHGTGFCFDSSFDVHKKCPLCELMFPPNYDQSKFEEHVESHWKV Sbjct: L55 DENVPTAPDPPSQHLRGHGTGFCFDSSFDVHKKCPLCELMFPPNYDQSKFEEHVESHWKV 734

Query: 757 CPMCSEQFPPDYDQQVFERHVQTHFDQNVLNFD 789
CPMCSEQFPPDYDQQVFERHVQTHFDQNVLNFD

Sbjct: 715 CPMCSEQFPPDYDQQVFERHVQTHFDQNVLNFD 747

Score = 104 (15.6 bits), Expect = 9.2e-02, Sum P(2) = 8.8e-02 Identities = 80/351 (22%), Positives = 157/351 (44%)

Query: 177 QLR---EQVGRMERELNH-EKERCDQLQAEQKGLTEVTQSLKMENEEFKKRFSDATSKAH 232 QLR EQV +E+ KE D + + +++ +++ENE+ KK+

40 Sbjct: 355 QLRKAEEQVQATRQEVVFLAKELSDAVNVRDRTMADL-HTARLENEKVKKQLADA---- 408

Query: 233 QLEEDIVSVTHKAIEKETE-

45 + + A++K+ + D+L+ +L++ E E L+ +L+ D
YK K +
Sbjct: 409 ----VAELKLNAMKKDQDKTDTLEHELRR---EVEDLKLRLQMAADH--YKEKFKECQ 457

50 Query: 292
IENTKLMSEVQTLKNLDGNKESVITHFKEEIGRLQLCLAEKENLQRTFLLTTSSKEDTCF 35L
+L ++ L + N +V T ++ G Q N T
T++S D
Sbjct: 458 ----RLQKQINKLSDQSANNNNVFT---KKTGNQQKVNDASVN--55 TDPATSASTVD--- 504

Query: 352 LKEQLRKAEEQVQ-ATRQEVVFLAKELSDAVNVRDRTMADLHTARLENEKVKKQLADAVA 430

+K AE T+ +V + KE++D ++ L + + K

+LA

Sbjct: 505

VKPSPSAAEADFDIVTKGQVCEMTKEIADKTEKYNKCKQLLQDEKAKCNKYADELAKMEL 564

5

10

Query: 411 ELKLNAMKKDQDKTDTLE----HELRREVED-LKLRLQMAAD--HYKEKFKECQ-RLQK 461

·RLQK 46』 - + K + E EL+R +E+ + +++ AD Y ++ +

R+

Sbjct: 565

KWKEQVKIAENVKLELAEVQDNYKELKRSLENPAERKMEDGADGAFYPDEIQRPPVRVPS 624

Query: 462 ---QINKLSDQSANNNNVFTKKTG--NQQKVNDASVNTDPATSASTVDVKPSPSAAEAD 515

15 + N + Q A N F++ G ++ D +V T P + +

+ ++

Sbjct: 625 WGLEDNVVCSQPARN---

FSRPDGLEDSEDSKEDENVPTAPDPPSQHLRGHGTGFCFDSS 681

20 Query: 516 FDIVTKGQVCEM 527

FD+ K +CE+

Sbjct: 682 FDVHKKCPLCEL 693

25 Pedant information for DKFZphamy2_2bl9, frame 2

Report for DKFZphamy2_2bl9.2

30 ELENGTHI 789

EMWI 90877-47

IpII 5.30

EHOMOLD TREMBL:HS338211_1 product: "tax1-binding protein

35 TXBP151"; Homo sapiens tax1-binding protein TXBP151 mRNA;

complete cds. D.D

EFUNCATO 99 unclassified proteins ES. cerevisiae, YOR216c3

EFUNCATI OB-O7 vesicular transport (golgi network, etc.)
ES-

40 cerevisiae, YDLO58wl 2e-13

EFUNCATI 30.03 organization of cytoplasm ES. cerevisiae YDL058wl 2e-13

45 EFUNCATI 30.04 organization of cytoskeleton ES. cerevisiae¬
YDR356wl 4e-l3

EFUNCATI 11.04 dna repair (direct repair, base excision repair

EFUNCATI D3.25 cytokinesis ES. cerevisiae, YHRD23w MY01 - myosin-l isoformI be-ll

55 [FUNCAT] OB.22 cytoskeleton-dependent transport [S. cerevisiae. YHRO23w MYO1 - myosin-l isoform] be-ll [FUNCAT] O3.04 budding. cell polarity and filament formation

EFUNCATE 1 genome replication, transcription, recombination and EM. jannaschii MJ13221 3e-08 repair EFUNCATD 98 classification not yet clear-cut ES. cerevisiae, YJR134c1 4e-08 EFUNCATI 03.19 recombination and dna repair ES. cerevisiae, YNL250wl 2e-07 EFUNCATI 03.13 meiosis ES. cerevisiae, YNL250wl 2e-07 EFUNCATD D3.Db cell growth ES. cerevisiae, YNLD79cD Ze-Db [FUNCAT] 03.07 pheromone response, mating-type determination, 10 sex-specific proteins IS. cerevisiae, YNLO79c3 2e-06 cerevisiae, YNLO79cl 2e-Db EFUNCATI D9-13 biogenesis of chromosome structure cerevisiae, YLRO86wl 5e-06 ES. cerevisiae, YPR141cl 2e-05 [[S. cerevisiae₁ YPR141cl 2e-05 EFUNCATE D3.22.DL cell cycle check point proteins cerevisiae, YGLO86w3 2e-05 EFUNCATD 30.05 organization of centrosome ES. cerevisiae. 20 YPR141c1 2e-05 EFUNCATI OB.16 extracellular transport ES- cerevisiaea YOR326wl le-04 EFUNCATD 09.25 vacuolar and lysosomal biogenesis cerevisiae, YOR326wl le-04 EFUNCATD 30.16 mitochondrial organization ES. cerevisiae. YALOLLwI Ze-04 EFUNCATD 06.07 protein modification (glycolsylation, acylation, myristylation, palmitylation, farnesylation and processing) ES. cerevisiae, YKL201c1 2e-04 30 EFUNCATI e amino acid metabolism and transport EM. genitalium. MG0423 4e-04 EFUNCATI 30.13 organization of chromosome structure cerevisiae, YDR285w3 7e-04 EFUNCATD n secretion and adhesion EM. jannaschii, MJD2913 35 0.001 EFUNCATI 05.04 translation (initiation, elongation and termination) [S. cerevisiae, YALO35w] [0.00] EBLOCKSD BLOO324D Tropomyosins proteins 40 EBLOCKSI PROD545E EBLOCKSI PRODO41F EZCOP] d2tmab_ 1.105.4.1.1 Tropomyosin Erabbit (Oryctolagus cuniculus) 5e-D5 CECI 3.6.1.32 Myosin ATPase 5e-16 45 **EPIRKU**J nucleus 2e-35 **EPIRKU** phosphotransferase 5e-10 duplication 2e-09 **EPIRKU EPIRKWI** citrulline 7e-09 [PIRKW] tandem repeat 2e-13 50 **EPIRKU** heterodimer 2e-D8 heart 2e-11 [PIRKW] endocytosis 3e-10 **EPIRKU** polymorphism le-09 **EPIRKU**J

EPIRKWI transmembrane protein be-12

55 EPIRKWI serine/threonine-specific protein kinase 5e-10
EPIRKWI cell wall 7e-09
EPIRKWI zinc finger 3e-10
EPIRKWI surface antigen be-08



```
EPIRKUI
                    DNA binding be-12
    [PIRKU]
                    metal binding 3e-10
                    muscle contraction 2e-13
    EPIRKUI
    EPIRKUI
                    brain 8e-D8
                    acetylated amino end 4e-09
    [PIRKW]
    [PIRKW]
                    actin binding 5e-16
    EPIRKU
                    endoplasmic reticulum 4e-09
    EPIRKUI
                    mitosis 3e-15
    [PIRKW]
                    microtubule binding 3e-15
10
    EPIRKUD
                    ATP 5e-16
    [PIRKW]
                    chromosomal protein 2e-08
    EPIRKUI
                    receptor 4e-10
    EPIRKUD
                    thick filament 2e-13
    [PIRKW]
                    phosphoprotein 5e-16
15
    EPIRKWI
                    glycoprotein 4e-10
                    skeletal muscle 7e-ll calcium binding 7e-09
    [PIRKW]
    [PIRKU]
    EPIRKWI
                    alternative splicing 3e-13
    EPIRKWI
                    DNA condensation 2e-08
20
    CPIRKW]
                    coiled coil 5e-1b
    [PIRKW]
                    P-loop 5e-16
                    heptad repeat 3e-13
    EPIRKWI
    EPIRKWD
                    methylated amino acid 2e-13
    [PIRKW]
                    basement membrane le-09
25
    EPIRKUI
                    immunoglobulin receptor 2e-09
    [PIRKW]
                   peripheral membrane protein 3e-10
    [PIRKW]
                   cardiac muscle 2e-11
    EPIRKU
                    extracellular matrix le-09
    EPIRKU
                   hydrolase 5e-16
30
    EPIRKWI
                   microtubule le-ll
    EPIRKWI
                   muscle le-09
                   membrane protein le-09
    [PIRKW]
    EPIRKUJ
                   EF hand 7e-09
    EPIRKWI
                   protein biosynthesis 4e-09
35
    EPIRKWD
                   cytoskeleton 3e-13
    [PIRKW]
                   hair 7e-09
    EPIRKWI
                   Golgi apparatus le-ll
    [PIRKW]
                   calmodulin binding 3e-10
    ESUPFAMD
              myosin heavy chain 5e-16
40
    EZUPFAMI
              conserved hypothetical Pll5 protein 4e-10
    ESUPFAMI
              IgA Fc receptor 7e-09
    ESUPFAM3
              centromere protein E 3e-15
              unassigned Ser/Thr or Tyr-specific protein kinases 5e-
    EZUPFAMJ
    70
45
    ESUPFAMD
              calmodulin repeat homology 7e-09
    [SUPFAM]
              myosin motor domain homology 5e-16
    ESUPFAMD
              alpha-actinin actin-binding domain homology 5e-10
    ESUPFAMD
              hypothetical protein MJO914 4e-O8
    [SUPFAM]
              tropomyosin be-09
50
    ESUPFAMJ
              plectin 5e-10
    ESUPFAMD
              trichohyalin 7e-09
    CSUPFAMI
              pleckstrin repeat homology le-O8
              ribosomal protein SLO homology 5e-10
    ESUPFAMI
    ESUPFAM3
              giantin 4e-13
55
    ESUPFAMI
              protein kinase homology 5e-10
    CSUPFAMD
              protein kinase C zinc-binding repeat homology le-08
    [SUPFAM]
              kinesin motor domain homology 3e-15
    ESUPFAM3
              human early endosome antigen 1 3e-10
```

	***	3 01/2042	,-									CI/IDVI	702030	
5	EKM] EKM] EKM] EZNЫ EZNЫ EZNЫ	FAMI FAMI	unass M5 pr cytos A11_/ L0W_(in MY(signed rotein skeled Alpha COMPLE ED_CO	d kin n 3e- tal k EXITY	esin 10 erat	in 4 3.3	e-07 0 %	prote	eins :	le-10			
10	SEQ SEG PRD COIL:	 cccee	eeec	cccee	eeeec	cccc	cccc	CCCC		CCCC	ccccc	cccee	FKVGWST	C
15	SEQ SEG PRD	RDYYT	FLWSF	PMPEHY	YVEGS	TVNC	VLAF	QGYYL	_PNDDG	EFYQ	FCYVTI	HKGEIF	RGASTPFQ	F •
20	COILS								_				· • • • • • • • • • • • • • • • • • • •	•
25	SEQ SEG PRD COILS	 hhhhh S	hhhhh	nhhhhh	hhhh	hhhh	hhhh	•xxx> hhhhh	(XXXXX hhhhhh	xxxx hhhhh	xxxxx; hhhhhhl	xxxxxx hhhhhh	KETAQLR (xx	h
	SEQ SEQ	QVGRM	ERELN	NHEKER	RCDQL	QAEQ	KGLT	EVTQS	SLKMEN	IEEFK	KRFSD.	ATSKAH		z
30	PRD COILS	7											CCCCCC	
35 ⁻	SEQ SEG PRD COILS	hhhhh	hhhhh	hhhhh	hhhh	hhhh	hhhh	hhhhh	hhhhh	hhhh	hhhhhl	nhhhhh	ENTKLMS	h
40	SEQ												 KEQLRKA	
	SEG PRD COILS	hhhhh	hhhhh	hhhhh	hhhh	hh hh	hhhhl	hhhhh	hhhhh	hhhhh	hhhhhl	nhhhhh	hhhhhhhh	h
45	SEQ SEG	EQVQA	TRQEV	VFLAK	ELSD	AVNV	RDRTI	MADLH	ITARLE	NEKVI	KQLA]	DAVAEL	KLNAMKK:	D •
50	PRD COILS											•		
55	SEQ SEG PRD COILS	hhhhhi	hhhhh	hhhhh	hhhhl	 nhhhi	hhhhl	hhhh	hhhhh	hhhhl	nhhhhl	nhhhhh	INNNVFTKI ihhhhhhhh	h
	SEQ	TGNQQI											DKTEKYNI	

WO 01/98454 PCT/IB01/02050 SEG PRD COILZ 5 CKQLLQDEKAKCNKYADELAKMELKWKEQVKIAENVKLELAEVQDNYKELKRSLENPAER SEQ SEG PRD րերի հերերի անագահաների անագահաների հերերի հերեր հերերի հերեր հերեր հերեր հերերի հերերի հերերի հերերի հերերի հերերի հերեր COILS 10 KMEGQNSQSPQCFKTCSEQNGYVLTLSNAQPVLQYGNPYASQETRDGADGAFYPDEIQRP SEQ SEG PRD 15 COILS PVRVPSWGLEDNVVCSQPARNFSRPDGLEDSEDSKEDENVPTAPDPPSQHLRGHGTGFCF SEQ SEG 20 PRD COILZ DSSFDVHKKCPLCELMFPPNYDQSKFEEHVESHWKVCPMCSEQFPPDYDQQVFERHVQTH SEQ 25 SEG PRD COILZ 30 SEQ FDQNVLNFD SEG PRD hcceeeccc COILS 35 (No Prosite data available for DKFZphamy2_2bl9.2) (No Pfam data available for DKFZphamy2_2bl9.2)

PCT/IB01/02050

DKFZphamy2_2c22

5 group: metabolism

DKFZphamy2_2c22 encodes a novel 364 amino acid protein with similarity to the 1-acyl-glycerol-3-phosphate acyltransferase of Zea mais.

It contains one leucine zipper. The protein is belived to play a role in fatty acid metabolism. It is ubiqitous expressed, with a slight predominance in uterus, placenta and foreskin.

- The new protein can find application in modulation of fatty acid metabolism and as a new enzyme for biotechnological production processes.
- weak similarity to 1-acyl-glycerol-3-phosphate acyltransferase
 (Zea
 mais)

perhaps complete cds.

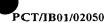
25 Sequenced by MediGenomix

Locus: /map="8"

30 Insert length: 3403 bp
Poly A stretch at pos. 3373, polyadenylation signal at pos. 3351

1 AGATGCTGCT GTCCCTGGTG CTCCACACGT ACTCCATGCG CTACCTGCTG 35 51-CCCAGCGTCG TGCTCCTGGG CACGGCGCCC ACCTACGTGT TGGCCTGGGG JOI GGTCTGGCGG CTGCTCTCCG CCTTCCTGCC CGCCCGCTTC TACCAAGCGC 151 TGGACGACCG GCTCTACTGC GTCTACCAGA GCATGGTGCT CTTCTTCTTC 201 GAGAATTACA CCGGGGTCCA GATATTGCTA TATGGAGATT TGCCAAAAA 251 TAAAGAAAT ATAATATT TAGCAAATCA TCAAAGCACA GTTGACTGGA BOL TTGTTGCTGA CATCTTGGCC ATCAGGCAGA ATGCGCTAGG ACATGTGCGC 40 351 TACGTGCTGA AAGAAGGGTT AAAATGGCTG CCATTGTATG GGTGTTACTT 401 TGCTCAGCAT GGAGGAATCT ATGTAAAGCG CAGTGCCAAA TTTAACGAGA 451 AAGAGATGCG AAACAAGTTG CAGAGCTACG TGGACGCAGG AACTCCAATG 501 TATCTTGTGA TTTTTCCAGA AGGTACAAGG TATAATCCAG AGCAAACAAA 45 551 AGTCCTTTCA GCTAGTCAGG CATTTGCTGC CCAACGTGGC CTTGCAGTAT **LOD TAAAACATGT GCTAACACCA CGAATAAAGG CAACTCACGT TGCTTTTGAT 651 TGCATGAAGA ATTATTTAGA TGCAATTTAT GATGTTACGG TGGTTTATGA** 701 AGGGAAAGAC GATGGAGGGC AGCGAAGAGA GTCACCGACC ATGACGGAAT 751 TTCTCTGCAA AGAATGTCCA AAAATTCATA TTCACATTGA TCGTATCGAC 50 BDL AAAAAAGATG TCCCAGAAGA ACAAGAACAT ATGAGAAGAT GGCTGCATGA B5D ACGTTTCGAA ATCAAAGATA AGATGCTTAT AGAATTTTAT GAGTCACCAG PDL ATCCAGAAAG AAGAAAAAGA TTTCCTGGGA AAAGTGTTAA TTCCAAATTA 951 AGTATCAAGA AGACTTTACC ATCAATGTTG ATCTTAAGTG GTTTGACTGC LODI AGGCATGCTT ATGACCGATG CTGGAAGGAA GCTGTATGTG AACACCTGGA 55 LOSI TATATGGAAC CCTACTTGGC TGCCTGTGGG TTACTATTAA AGCATAGACA 1101 AGTAGCTGTC TCCAGACAGT GGGATGTGCT ACATTGTCTA TTTTTGGCGG 1151 CTGCACATGA CATCAAATTG TTTCCTGAAT TTATTAAGGA GTGTAAATAA 1201 AGCCTTGTTG ATTGAAGATT GGATAATAGA ATTTGTGACG AAAGCTGATA

WO 01/98454



		01/20454				1 C 1/1B01/02030
	1251	TGCAATGGTC	TTGGGCAAAC	ATACCTGGTT	GTACAACTTT	AGCATCGGGG
	1301	CTGCTGGAAG	GGTAAAAGCT	AAATGGAGTT	TCTCCTGCTC	TGTCCATTTC
	1351	CTATGAACTA	ATGACAACTT	GAGAAGGCTG	GGAGGATTGT	GTATTTTGCA
	1401	AGTCAGATGG	CTGCATTTTT	GAGCATTAAT	TTGCAGCGTA	TTTCACTTTT
5	1451	TCTGTTATTT	TCAATTTATT	ACAACTTGAC	AGCTCCAAGC	TCTTATTACT
	1501	AAAGTATTTA	GTATCTTGCA	GCTAGTTAAT	ATTTCATCTT	TTGCTTATTT
	1551	CTACAAGTCA	GTGAAATAAA	TTGTATTTAG	GAAGTGTCAG	GATGTTCAAA
	1601	GGAAAGGGTA	AAAAGTGTTC	ATGGGGAAAA	AGCTCTGTTT	AGCACATGAT
	1651	TTTATTGTAT	TGCGTTATTA	GCTGATTTTA	CTCATTTTAT	ATTTGCAAAA
10	1701	TAAATTTCTA	ATATTTATTG	AAATTGCTTA	ATTTGCACAC	CCTGTACACA
	1751	CAGAAAATGG	TATAAAATAT	GAGAACGAAG	TTTAAAATTG	TGACTCTGAT
	1801	TCATTATAGC	AGAACTTTAA	ATTTCCCAGC	TTTTTGAAGA	TTTAAGCTAC
	1851	GCTATTAGTA	CTTCCCTTTG	TCTGTGCCAT	AAGTGCTTGA	AAACGTTAAG
	1901	GTTTTCTGTT	TTGTTTTGTT	TTTTTAATAT	CAAAAGAGTC	GGTGTGAACC
15	1951	TTGGTTGGAC	CCCAAGTTCA	CAAGATTTTT	AAGGTGATGA	GAGCCTGCAG
15	5007	ACATTCTGCC	TAGATTTACT	AGCGTGTGCC	TTTTGCCTGC	TTCTCTTTGA
	2051	TTTCACAGAA	TATTCATTCA	GAAGTCGCGT	TTCTGTAGTG	TGGTGGATTC
	5707	CCACTGGGCT	CTGGTCCTTC	CCTTGGATCC	CGTCAGTGGT	GCTGCTCAGC
	2151	GGCTTGCACG	CAGACTTGCT	AGGAAGAAAT	GCAGAGCCAG	CCTGTGCTGC
20	5507	CCACTTTCAG	AGTTGAACTC	TTTAAGCCCT	TGTGAGTGGG	CTTCACCAGC
20	2251	TACTGCAGAG	GCATTTTGCA	TTTGTCTGTG	TCAAGAAGTT	CACCTTCTCA
	5307	AGCCAGTGAA	ATACAGACTT	AATTTGTCAT	GACTGAACGA	ATTTGTTTAT
	2351	TTCCCATTAG	GTTTAGTGGA	GCTACACATT	AATATGTATC	GCCTTAGAGC
	2401	AAGAGCTGTG	TTCCAGGAAC	CAGATCACGA	TTTTTAGCCA	TGGAACAATA
25	2451	TATCCCATGG	GAGAAGACCT	TTCAGTGTGA	ACTGTTCTAT	TTTTGTGTTA
	2501	TAATTTAAAC	TTCGATTTCC	TCATAGTCCT	TTAAGTTGAC	ATTTCTGCTT
	2551	ACTGCTACTG	GATTTTTGCT	GCAGAAATAT	ATCAGTGGCC	CACATTAAAC
	5607	ATACCAGTTG	GATCATGATA	AGCAAAATGA	AAGAAATAAT	GATTAAGGGA
	2651	AAATTAAGTG	ACTGTGTTAC	ACTGCTTCTC	CCATGCCAGA	GAATAAACTC
30	2701	TTTCAAGCAT	CATCTTTGAA	GAGTCGTGTG	GTGTGAATTG	GTTTGTGTAC
	2751	ATTAGAATGT	ATGCACACAT	CCATGGACAC	TCAGGATATA	GTTGGCCTAA
	5907	TAATCGGGGC	ATGGGTAAAA	CTTATGAAAA	TTTCCTCATG	CTGAATTGTA
	2851	ATTTTCTCTT	ACCTGTAAAG	TAAAATTTAG	ATCAATTCCA	TGTCTTTGTT
	2901	AAGTACAGGG	ATTTAATATA	TTTTGAATAT	AATGGGTATG	TTCTAAATTT
35	2951	GAACTTTGAG	AGGCAATACT	GTTGGAATTA	TGTGGATTCT	AACTCATTTT
	3001	AACAAGGTAG	CCTGACCTGC	ATAAGATCAC	TTGAATGTTA	GGTTTCATAG
	3051	AACTATACTA	ATCTTCTCAC	AAAAGGTCTA	TAAAATACAG	TCGTTGAAAA
	3707	AAATTTTGTA	TCAAAATGTT	TGGAAAATTA	GAAGCTTCTC	CTTAACCTGT
	3151	ATTGATACTG	ACTTGAATTA	TTTTCTAAAA	TTAAGAGCCG	TATACCTACC
40	3507	TGTAAGTCTT	TTCACATATC	ATTTAAACTT	TTGTTTGTAT	TATTACTGAT
	3251	TTACAGCTTA	GTTATTAATT	TTTCTTTATA	AGAATGCCGT	CGATGTGCAT
	3307	GCTTTTATGT	TTTTCAGAAA	AGGGTGTGTT	TGGATGAAAG	TAAAAAAAAA
	3351	AAATAAAATC	TTTCACTGTC	TCTAAAAAAA	AAAGAAAAA	AAAAAAAAA
	3401	AAA				
45						

BLAST Results

50 No BLAST result

Medline entries

No Medline entry

PRD



Peptide information for frame 3

5 ORF from 3 bp to 1094 bp; peptide length: 364 Category: similarity to known protein Classification: Metabolism Prosite motifs: LEUCINE_ZIPPER (105-126) 10 I MLLSLVLHTY SMRYLLPSVV LLGTAPTYVL AWGVWRLLSA FLPARFYQAL 51 DDRLYCVYQS MVLFFFENYT GVQILLYGDL PKNKENIIYL ANHQSTVDWI JOJ VADILAIRAN ALGHVRYVLK EGLKWLPLYG CYFARHGGIY VKRSAKFNEK 151 EMRNKLQSYV DAGTPMYLVI FPEGTRYNPE QTKVLSASQA FAAQRGLAVL 201 KHVLTPRIKA THVAFDCMKN YLDAIYDVTV VYEGKDDGGQ RRESPTMTEF 15 251 LCKECPKIHI HIDRIDKKDV PEEQEHMRRW LHERFEIKDK MLIEFYESPD 301 PERRKRFPGK SVNSKLSIKK TLPSMLILSG LTAGMLMTDA GRKLYVNTWI 351 YGTLLGCLWV TIKA 20 BLASTP hits No BLASTP hits available 25 Alert BLASTP hits for DKFZphamy2_2c22, frame 3 No Alert BLASTP hits found 30 Pedant information for DKFZphamy2_2c22, frame 3 Report for DKFZphamy2_2c22.3 . 35 ELENGTHD 364 42072-47 EMMI [pl] 9-18 EHOMOLI TREMBL: CEAF3136_1 gene: "F28B3.5"; Caenorhabditis 40 elegans cosmid F28B3. 2e-36 EFUNCATI 99 unclassified proteins ES. cerevisiae, YDROLAcl 7e-13 EFUNCATI D1.06.01 lipid, fatty-acid and sterol biosynthesis [S. cerevisiae, YDL052c] 4e-05 45 **EFUNCATI** 30.99 other cellular organization ES. cerevisiae, YDL052cJ 4e-05 **EBLOCKZD** BLD15P3V EBLOCKSI BPDD989A **EPIRKWI** transmembrane protein 2e-11 50 ESUPFAMI probable membrane protein YBR042c 2e-11 EPROSITEJ LEUCINE_ZIPPER 1 Alpha_Beta LOW_COMPLEXITY EKWI 3.57 % 55 SEQ MLLSLVLHTYSMRYLLPSVVLLGTAPTYVLAUGVURLLSAFLPARFYQALDDRLYCVYQS SEG





5	SEQ SEG PRD	MVLFFFENYTGV@ILLYGDLPKNKENIIYLANH@STVDWIVADILAIR@NALGHVRYVLK
J	SEQ SEG PRD	EGLKWLPLYGCYFAQHGGIYVKRSAKFNEKEMRNKLQSYVDAGTPMYLVIFPEGTRYNPE hhhcccccceeeccceeeeeccccchhhhhhhhhhhhh
10	SEQ SEG PRD	QTKVLSASQAFAAQRGLAVLKHVLTPRIKATHVAFDCMKNYLDAIYDVTVVYEGKDDGGQ
15	SEQ SEG PRD	RRESPTMTEFLCKECPKIHIHIDRIDKKDVPEE@EHMRRWLHERFEIKDKMLIEFYESPD
20	SEQ SEG PRD	PERRKRFPGKSVNSKLSIKKTLPSMLILSGLTAGMLMTDAGRKLYVNTWIYGTLLGCLWV cccccccccchhhhhhhhhhhhhhhhhhhhhhhhhhh
25	SEQ SEG PRD	TIKA hccc

Prosite for DKFZphamy2_2c22.3

105->127 LEUCINE_ZIPPER PD0C00029 30 PS00029

(No Pfam data available for DKFZphamy2_2c22.3)

DKFZphamy2_2fl8

5 group: signal transduction

DKFZphamy2_2fl& encodes a novel 215 amino acid protein with similarity to sodium channel protein betal of Rattus norvegicus.

The sodium channel protein beta 1 of Rattus norvegicus is crucial in the assembly expression and functional modulation of the heterotrimeric complex of the rat brain sodium channel. The expression of the new protein seems to be restricted to brain all matching ESTs isolated so far derive from there.

The new protein can find application in modulating the sodium channel beta 1, studying the expression profile in neurodegenerative diseases and of amygdala -specific genes.

20

15

similarity to sodium channel protein betal (Rattus norvegicus)

L CAGGGCTGAC AGCACACG GCCTGGGGGC CTAGAGAAGG ATTGCTGATC

Pedant: SIGNAL_PEPTIDE

25

Sequenced by MediGenomix

Locus: unknown

30 Insert length: 4052 bp
Poly A stretch at pos. 4035, no polyadenylation signal found

51 ACCTGCCACC CAGGGTCGGG GCCCCGCACC ATCCGGGGGC GAGCTCCCGG 35 101 GAAGGGGCTC CCCCTCTACA CCCACCCCC AACCTCTGAC ATCGCCGGCC 151 GAACGGGAGC TGCCGCTTCC TTCCCGGCCC CGCTGCACCT CCCCAGGGAG 201 CCGAGGGCGG GCGTGGACGG GACCGACGTG GAACGCATTC TGTAGCCCAG 251 ACGGGCGGCC CCGGCGGCTT CGGGAGTGGG GTCACGCCCA GCTGGAGAAG 3D1 CAGTTAGGGC GGACGAAGCA GGAGCCGCGG GGCTGGGAGG ATTCCAGTCG 40 351 GAACGCAACC GATCCTGGGG AGGCGAGAGG TGAATCAACC TGGACCCTTC 4D1 CACAGCCTGG CTGCTAGGCC AGCAGTGCGA CTCCCTTCCG AGCTGAGCTT 451 ACCCTGGGCG CAAACGAGCG AGGCAGGGGC GCGAGTGGAA GCTGGAGTTC 5D1 CGGGGTGGGC GGGGAGGCGA CTGTCCGTGG TGCTGAGCGC CGGCGAGAGC 551 GGGCGCGAG CGGCTGATCA GCTCCCTCGA ACTGGGGAGG TCCAGTGGGG 45 LOL TCGCTTAGGG CCCAAAGCCC CCGCCCGGCT CCAAAAGCTC CCAGGGCCTC LSL CCCAGGCACC GGTGCTCGGC CCTTCCTTCG GTCAGAAAGT CGCCCCCTGG 7D1 GGGCAGTTCG TCCCAAAGGG TTTCCTCGAA AGAATCTGAG AGGGCGCAGT 701 GGGCAGTTCG TCCCAAAGGG TTTCCTCGAA AGAATCTGAG AGGGCGCAGT
751 CCTTGACCGA GGGAATCTCT CTGTGTAGCC TTGGAAGCCG CCAGCCCCAG
BD1 AAGATGCCTG CCTTCAATAG ATTGTTTCCC CTGGCTTCTC TCGTGCTTAT
B51 CTACTGGGTC AGTGTCTGCT TCCCTGTGTG TGTGGAAGTG CCCTCGGAGA
901 CGGAGGCCGT GCAGGGCAAC CCCATGAAGC TGCGCTGCAT CTCCTGCATG
951 AAGAGAGAG AGGTGGAGGC CACCACGGTG GTGGAATGGT TCTACAGGCC
1001 CGAGGGCGGT AAAGATTTCC TTATTTACGA GTATCGGAAT GGCCACCAGG
1051 AGGTGGAGAG CCCCTTTCAG GGGCGCCTGC AGTGGAATGG CAGCAAGGAC
1101 CTGCAGGACG TGTCCATCAC TGTGCTCAAC GTCACTCTGA ACGACTCTGG
1151 CCTCTACACC TGCAATGTGT CCCGGGAGTT TGAGTTTGAG GCGCATCGGC
1201 CCTTTGTGAA GACGACGCGG CTGATCCCCC TAAGAGTCAC CGAGGAGGCT 50 55



	1301 1301		TCACCTCTGT ACCTTGTGGC		ATCATGATGT GATGATATAT	ACATCCTTCT TGCTACAGAA
	1351	AGGTCTCAAA		GCAGCCCAAG	AAAACGCGTC	TGACTACCTT
	1401	GCCATCCCAT	CTGAGAACAA	GGAGAACTCT	GCGGTACCAG	TGGAGGAATA
5	1451	GAACAGGAGC		GAGGTGGCCT	GAACACCTGA	GGGACTGGAC
	1501	ATCCCATGTT	CAGCAATGTC	AATGGCATCA	GGAGGGCGCC	CCAAGGGCCC
	1551	CATCGCTTCC	_	CCATTGTTCT	GTTCATTCAT	TCATCCATAC
	1201	ATCCACCTGC		TCACCTCTGA		CCATCAGACC
	1651	TCTACGCACC		GCCAGAACTG	AGAAGCCAAC	ATTTCTACAT
10	1701	AGACTCAACC		CTAGTTTTCC	AACAAGACAC	TCCAAAGCCA
	1751	ACTGGATTTC	TCCCCTGTGC	TCCAAATGAC	TTTGTACAAG	TGCTGGAGTT
	1801	AGCACCTCCC		ACTGGCTGGA	ACTGGTTCAT	TCTCCATTAC
	1851	TGCAAGAGAA	TGGAAGTCTT	AATAGAAGGA	AGCAGGAGTG	ATTAGTTCGG
	1901	GTTAAAGCAA	AAGTGTGTCA	TGAACTTGGA	TTCCCTGAAG	TCAGTTTTGT
15	1951	CAGGTTCATG		CTACAGCATC	AGAGTGAAGC	ACGCCTGTCT
	5007	AGGTTCTCCA	·· · · · · -	GATCCTGAAG		CATGCTCTCT
	2051	GGAGCTTAGT	ACTCCAGAGC	TAGATCCTGA	TGGGTCTCTA	AGGTTCCCTC
	5707	CAAGAAGACA	AGGACAGGAG	ACTTGGGAAG	GACCAATGGT	AATTTAAGTG
	2151	GCTCTTAAAA	AGTCATGCAA	CATGTTTCTG	GACACGTTCC	TGATCCTATT
20	5501	GCGATAATGT	ATGTGTGCCC	TCCCTGTGGG		GGGCATTAGG
	2251	ACTGAAATTC	CTGAGTTCTT	CCTCTCAAAA	TTTCTGTGCA	CCAGTATTAT
	5301	TCCTCATTTT	ACATACAGGA	GGCAACTAAG	ACTCATACAG	GGCTCAACTG
	2351	AATAAGAGGC	TTAAGAGGAT	AAACTGGAGC	AGAAATAAGC	CTTAGGTGCT
	2401	GCCCAGTTTA	CACTTCCTGG	GATGGATGTT	TTTGTTTGTT	TTGTTTTTTG
25	2451	TTTTTTTTGT	TTGAGATGGA	GTCTCACTCT	GTCACCTAGG	CTAGAGTGCA
	2501	GTGGTGTGAT	CTCGGCTCAC	TGCAACCTCT	GCCTCTTGGG	TTCAAGCAAT
	2551	TCTCATGCCT	CGGCCTCTCC	AGTAGCTGGG	ATTACAGGTG	TGCACCACCA
	5201	CGCCTGGCTA	AATTTTGTAT	TTTTAGTACA	GACAGGGTTT	GACTATGTTG
	2651	GCCAGGCTAG	TCTTGAACTC	CTGACCTCAA	ATGACCCACC	CACCTCAGCC
30	2701	TCCCAAAGTG	CTGAGATTAC	AGGCGTGAGG	CACTGCGCCC	GGTGGATAAC
	2751	TTTGTTTCTG	AAAAGACTGA	CATTGAACTT	GTCTATGGCA	ATGCTTCTTT
	5901	CACAAGCACG	GACTGGGCTG	AGGTCAACTC	TGATAGATTC	AGATGACTAG
	2851	AAATTGGCCA	AAAAAGCAGG	GAGAAGAACA	TGAGGTAGAC	TTAAAGAACT
	5407	TCCTTTATGT	AAAGATCTGT	GACTCTGAAA	TATCCTCCAA	AAGGAGAGTG
35	2951	CATCTGAGAC	TGATATTTAA	ACTAAGAAAA	ATGTTTAGTC	TGAGATGGAT
	3007	CATAAGTAAA	TGAGCAGTGT	GAGAGGGGAG	GGATGGGTAG	GTGCTTTCCA
	3051	AATACTTCGC	CTATGAATGC	ATAATTTTCA	GATTTTTTC	CCCTAGATTT
	3707	TGAGGGAGCA	GAGAAACTGG	AAAAAACTTT	AGTCAATATC	TCGTGTTTCA
40	3151	TTTTAATTAA	GTGACAGGTC	CAAGTGTGAC	ATCCTTCAGC	ACCCAGGGAC
40	3507	AAGAGAGGG	AAAGATGCTT	TATGGAATGT	AAGAAGATGA	AGGTGACTGG
	3251	GATTCAGCGA	GAGAGAGGTC	CCTCAGACCT	GGGACCTCCC	TTTATAGGGA
	3301	AAGACCATAT	TCCATAGGTT	TAGGGCTTTA	CCTTAAAAGC	TCATTTTTT
	3351	CATTCTTCCA	TCCCTAGGAA	AGTACTTAAA	ACCAGACTTT	TAAATTTTTA
AE	3401	TTTATTTATT	ATTATTTTT	TGAGACAGAT	TCTCACTCTG	TCTCCCAGGC
45	3451	TAGAGTGCAG	TGGTGCAATC	TCAGCTCACT	GCAGCCTCAA	CTGCCCCAGG
	3501	TTTAAGCAAT	CCTCCCACCT	CAGCCCCCAG	GTAACTGGGA	
	3551	GCACCACCAT TGCCATGTCG	GCCTGGCTAA	TTTTTGTATT	TTATGTAGAG	
	3607		CCCAGGCTGA	TCTTGAACTC	CTGGGCTCAA	GCAATCTGCC
50		AGCCTCAGCC	TCTCAAAGTG	CTGGGATTAC	AGGCCTGAGC	AACTGTGCCT
50		GGCCCAAAAC	CAGACCGTTA	ACACATTAAA	GAGTCTGATT	TTGTTGAAGA
	3751	AAATATTTGC	AATAAATTCA	AGACTCTTCT	TATTGGTAAT	TTTCCACACA
	3801	ATCCCTCTGA	AATAAGGGAG	AGGATATAGA	CCTTTTTAAC	TTTATAGTTA
	3851	GAAAAATTGG	CCTCAGTGTG	AAATTTTTCC	AGTCCCATAG	CTCATGGATG
55			GCGGTAGTAG	CAAGATGCTT	ACTACCACAC	
55	3951		TAGCTCGTGT	ATCTAAGTTG	AACCCGGCAG	TATGCATGAT
	4001 4051	TGCCTTTTTC	TCTTCTTTTT	AAAAAACCC	AACTCAAAAA	AAAAAAAAA
	TCUF	nn				



BLAST Results ------

No BLAST result 5

Medline entries ______

10

Isom LL, De Jongh KS, Patton DE, Reber BF, Offord J, Charbonneau

Walsh K, Goldin AL, Catterall

WA.; Primary structure and functional expression of the beta l 15 subunit

the rat brain sodium channel. Science 1992 May 8:256(5058):839-42

96235151: 20

Belcher SM, Howe JR.; Cloning of the cDNA encoding the sodium

beta I subunit from rabbit. Gene 1996 May 8:170(2):285-6

25

McClatchey AI, Cannon SC, Slaugenhaupt SA, Gusella JF.; The cloning and

expression of a sodium channel beta

1-subunit cDNA from human brain. Hum Mol Genet 1993 Juni2(6):745-

30

35

45

Peptide information for frame 3

ORF from 804 bp to 1448 bp; peptide length: 215 Category: similarity to known protein

Classification: Transmembrane proteins unclassified 40

I MPAFNRLFPL ASLVLIYWVS VCFPVCVEVP SETEAV@GNP MKLRCISCMK

THE TOTAL TENDESCRIPTION STREET AND STREET A

151 EDFTSVVSEI MMYILLVFLT LWLLIEMIYC YRKVSKAEEA AQENASDYLA

507 Ibzenkenza Abaee

BLASTP hits 50

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_2fl8, frame 3

PIR:JC4788 sodium channel protein betal chain - rabbit, N = l, 55 Score = 434, P = 8.3e-41

```
PIR:A55734 sodium channel, voltage-gated, beta-1 chain precursor
    human, N = 1, Score = 428, P = 3.6e-40
 5
    PIR:A42737 sodium channel beta 1 subunit - rat, N = 1, Score =
    429, P =
    2.8e-40
10
    >PIR:JC4788 sodium channel protein betal chain - rabbit
                Length = 218
      HSPs:
15
     Score = 434 (65.1 bits), Expect = 8.3e-41, P = 8.3e-41
      Identities = 100/214 (46%), Positives = 129/214 (60%)
              ЪП
. 20
    LASLVLIYWVSVCFPVCVEVPSETEAVQGNPMKLRCISCMKREEVEATTVVEWFYRPEGG 69
                 LA +V VS + CVEV SETEAV G K+ CISC +R E A T
    EW +R +6
    Sbict:
    LAFVVGAALVSSAWGGCVEVDSETEAVYGMTFKILCISCKRRSETTAETFTEWTFRQKGT 64
25
             70 KDFL-IYEYRNGHQEVESP--FQGRLQUNGS---
    Query:
   KDLQDVSITVLNVTLNDSGLYTCNVS 123
                 ++F+ I Y N ++E
                                     F+GR+ WNGS KDLQD+SI + NVT N
    ZG Y C+V
30
    Sbjct:
              Ь5
    EEFVKILRYENEVLQLEEDERFEGRVVWNGSRGTKDLQDLSIFITNVTYNHSGDYQCHVY 124
    REFEFEAHRPFVKTTRLIPLTVTEEAGEDFTZVVZEIMMYIXXXXXXXXXIEMIYCYRK 1A3
                 R FE + - + I L V ++ A D S+VSEIMMY+
35
    EM+YCY+K
    Sbict:
    RLLSFENYEHNTSVVKKIHLEVVDKANRDMASIVSEIMMYVLIVVLTIWLVAEMVYCYKK 184
    Query: 184 VSKAEEAA-QENASDYLAIPSENKEN-SAVPVEE 215
40
                 ++ A EAA QENAS+YLAI SE+KEN + V V E
             185 IAAATEAAAQENASEYLAITSESKENCTGVQVAE 218
    Sbjct:
45
                Pedant information for DKFZphamy2_2flb, frame 3
                Report for DKFZphamy2_2fl8.3
50
    ELENGTHD 215
             24702.40
    EMMI
              4.69
    [[pI]
    EHOMOL
                  PIR: JC4788 sodium channel protein betal chain -
55
    rabbit 3e-41
    CBLOCKSI BLOO401D Prokaryotic sulfate-binding proteins
    EBFOCKZD Bb00240
```

PCT/IB01/02050 WO 01/98454 2.1.1.1.1 Myelin membrane adhesion dlneu__ LZCOD1 molecule PO [ra 2e-43 [PIRKW] Schwann cell 2e-07 **EPIRKW** transmembrane protein le-40 myelin 2e-07 **EPIRKWI** phosphoprotein 5e-07 [PIRKW] glycoprotein le-40 [PIRKW] structural protein 2e-07 **EPIRKW** muscle le-40 [PIRKU] 10 **EPIRKUD** membrane protein 5e-07 immunoglobulin homology 2e-D7 ESUPFAMI **ESUPFAMD** myelin PO protein 2e-07 IG (immunoglobulin) superfamily EPFAM1 EKWI All_Beta 15 EKW] 3D EKW3 SIGNAL_PEPTIDE 23 EKWI LOW_COMPLEXITY 4-65 % SEQ MPAFNRLFPLASLVLIYWVSVCFPVCVEVPSETEAVQGNPMKLRCISCMKREEVEATTVV 20 lneu------CEEEECCEEETTTbCEEECE-EEECCCCCCCCEE 25 SEQ EWFYRPEGGKDFLIYEYRNGHQEVESPFQGRLQWNGSKDLQDVSITVLNVTLNDSGLYTC SEG lneu-EEEEEETTTCCCEEEEEETTEEEETTTTTTTEEECCBGGCBCCEEECCbTTTTTEEEEE NVSREFEFEAHRPFVKTTRLIPLRVTEEAGEDFTSVVSEIMMYILLVFLTLWLLIEMIYC 30 SEQ SEG -----XXXXXXXXXXXX lneu-SEG YRKVSKAEEAAGENASDYLAIPSENKENSAVPVEE 35 SEG Ineu-40 (No Prosite data available for DKFZphamy2_2fl8.3) Pfam for DKFZphamy2_2fl8.3 45

55

50

HMM

Y+C+V Query YRNG ++ E+ ++ R++++G ++ +++T+ +++ +DSG

HMM_NAME IG (immunoglobulin) superfamily

. 77 YRNGHQEV--ESPFQGRLQWNGSKDLQDVSITVLNVTLNDSGLYTCNV 122

yrNgqpipssegyWytRweqqgRYsisifqLtIisWepeDsGtYWCmV

WO 01/98454

,

PCT/IB01/02050

DKFZphamy2_2f22

5 group: nucleic acid management

DKFZphamy2_2f22 encodes a novel 479 amino acid protein with similarity to YDL153c of Saccharomyces cerevisia.

10 The novel protein is ubiquitously expressed. YDL153c is involved in transcriptional silencing.

The new protein can find application in modulation of transcription, e.g. transcriptional silencing.

15

putative protein

probably complete cds.
20 perhaps differential polyadenylation
YDL153c is involved in transcriptional silencing

Sequenced by MediGenomix

25 Locus: /map="4"

Insert length: 2019 bp

Poly A stretch at pos. 2000, polyadenylation signal at pos. 1981

30 L GGAGTCTGCA AACTCCGGTG GTAGGGGAGC GCGCTGCTGT TTAGAGCCAC 51 GAGTTACCGG AGCGCCTGAT TCCTGCGCCG AAGTCAGTGG TGGCCGAAAG JOD TCCGGAGTCG CTGTAAAACC TGAGATTGTG AGCCATGGTG GGGAGATCCC 151 GGCGGCGCG AGCAGCTAAG TGGGCAGCTG TGCGAGCCAA GGCAGGTCCC 35 201 ACGCTCACCG ACGAAAATGG AGATGATTTA GGATTGCCAC CCTCACCAGG 251 GGACACCAGC TACTACCAAG ATCAGGTAGA TGACTTTCAT GAGGCACGAT BDL CCCGGGCCGC CTTAGCTAAG GGCTGGAATG AAGTACAGAG TGGAGACGAG 351 GAGGATGGCG AGGAGGAGGA GGAGGAGGTG CTAGCCCTAG ATATGGACGA 401 TGAGGACGAC GAAGATGGAG GGAATGCGGG GGAGGAGGAG GAGGAGGAGA 45% ATGCCGATGA TGATGGTGGG AGCTCCGTGC AAAGTGAAGC TGAGGCCTCT 40 501 GTGGATCCCA GTTTGTCGTG GGGTCAGAGG AAAAAACTTT ACTATGACAC 551 GGACTATGGT TCCAAGTCCC GAGGCCGGCA GAGTCAACAG GAGGCAGAGG LOT AGGAGGAAAG AGAGGAGGAG GAGGAGGCAC AGATCATTCA GCGGCGCCTA **L51 GCCCAAGCGC TGCAAGAGGA TGATTTTGGT GTCGCCTGGG TTGAGGCCTT** 45 701 TGCAAAACCA GTGCCTCAGG TAGATGAGGC TGAGACACGG GTCGTGAAGG 751 ATTTGGCTAA AGTTTCAGTG AAAGAGAAGC TGAAAATGTT GCGAAAGGAA BOD TCACCAGAAC TCTTGGAGCT GATAGAAGAC CTGAAAGTCA AGTTGACAGA B51 GGTTAAGGAT GAGCTGGAGC CATTGTTAGA GTTGGTGGAA CAAGGGATCA 9D1 TTCCACCCGG AAAAGGAAGC CAATACTTGA GGACCAAGTA CAACCTCTAC 951 TTGAATTATT GCTCGAACAT CAGTTTTTAT TTGATCCTGA AAGCTAGGAG 50 LDDL AGTCCCAGCA CATGGACATC CTGTCATAGA AAGGCTTGTT ACCTACCGAA 1051 ATTTGATCAA CAAGCTGTCC GTTGTGGATC AGAAGCTGTC CTCAGAAATT 1101 CGTCATCTGT TGACACTTAA GGATGATGCT GTAAAGAAAG AACTGATTCC 1151 AAAAGCAAAA TCCACCAAGC CCAAACCAAA GTCTGTTTCA AAGACTTCTG 1201 CTGCTGCCTG TGCTGTTACA GATCTTTCTG ATGATTCTGA TTTTGATGAA 55 1251 AAAGCAAAAC TGAAGTACTA TAAAGAAATA GAAGACAGGC AAAAGCTAAA 1301 GAGAAAGAAA GAAGAAAATA GCACTGAAGA ACAGGCTCTT GAAGATCAAA 1351 ATGCAAAGAG AGCTATTACC TATCAAATTG CTAAAAATAG GGGACTTACT

LHOL CCTAGGAGAA AGAAGATTGA TCGCAATCCC AGAGTGAAAC ACAGAGAGAA
LH5L GTTCAGAAGA GCCAAAATTA GAAGAAGAGG CCAGGTTCGT GAAGTTCGTA
L5DL AAGAAGAGCA ACGTTATAGT GGTGAATTAT CTGGCATTCG TGCAGGAGTT
L55L AAAAAGAGCA TTAAGCTTAA ATGAAGTTTT TGCTTAGCAT AAGGTTTTTG
LLDL GCAGTTTTGG ATCAATAAAT TTTTACTTTT AACTAAAGTC ATTGTATTAA
LL5L TATATAATAC TTTAAATTTT AAAAATTCTT GTCCACAAGG AAATTTGTCT
L7DL GGGTTATTGG ACAATTTATA AGAACTATGG GAGCAATATG AAGGTGCTTG
L75L AGAAAAGAGA TGATGTTGAA GTTTTCCAAT ATTCTGTTGA AGTTTTCCAA
LBDL TATTAAGTAT TAGCTTAGGG AAATTTCACA GTTCATTGTG GAGTGTTAAA
LBSL CTTAGAACAT GTGTAACTTT TCACATAAAG AGAATGCATC TTTGACAGTT
L9DL ATCTTATTTG TAAGGCAGCC TATAAAAATAG TTCTGAAGTA TTTTATTTAC
L95L CTAACTATAA TTATTGGGCC AGATACTTGT TAATAAATGG GCTTAATGTC

15

BLAST Results

No BLAST result

20

Medline entries

25 No Medline entry

Peptide information for frame 3

30

ORF from 135 bp to 1571 bp; peptide length: 479 Category: similarity to unknown protein Classification: Nucleic acid management

35

40

45

I MVGRSRRRGÅ AKWAAVRAKA GPTLTDENGD DLGLPPSPGD TSYYQDQVDD 51 FHEARSRAAL AKGWNEVQSG DEEDGEEEEE EVLALDMDDE DDEDGGNAGE 101 EEEEENADDD GGSSVQSEAE ASVDPSLSWG QRKKLYYDTD YGSKSRGRQS 151 QQEAEEEERE EEEEAQIIQR RLAQALQEDD FGVAWVEAFA KPVPQVDEAE 201 TRVVKDLAKV SVKEKLKMLR KESPELLELI EDLKVKLTEV KDELEPLLEL 251 VEQGIIPPGK GSQYLRTKYN LYLNYCSNIS FYLILKARRV PAHGHPVIER 301 LVTYRNLINK LSVVDQKLSS EIRHLLTLKD DAVKKELIPK AKSTKPKPKS 351 VSKTSAAACA VTDLSDDSDF DEKAKLKYYK EIEDRQKLKR KKEENSTEEQ 401 ALEDQNAKRA ITYQIAKNRG LTPRRKKIDR NPRVKHREKF RRAKIRRRGQ 451 VREVRKEEQR YSGELSGIRA GVKKSIKLK

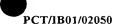
BLASTP hits

50

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_2f22, frame 3

55 PIR:S67701 hypothetical protein YDL153c - yeast (Saccharomyces cerevisiae), N = 4, Score = 134, P = 1.8e-08



PIR:TO8694 hypothetical protein DKFZp5640092.1 - human (fragment), N = 1.2 Score = 141.7 P = 5.8e-07

5 TREMBL:SPBC3BA_9 gene: "SPBC3BA.09"; product: "hypothetical protein"; S.pombe chromosome II cosmid c3BA., N = 2, Score = 164, P = 6.2e-13

>TREMBL:SPBC3BA_9 gene: "SPBC3BA-09"; product: "hypothetical
protein";

S-pombe chromosome II cosmid c3B8-Length = 597

15 HSPs:

Score = 164 (24.6 bits), Expect = 6.2e-13, Sum P(2) = 6.2e-13 Identities = 44/126 (34%), Positives = 68/126 (53%)

20

Query: 367 DSDFDEKAKLKYYKEIEDRQKLKRK-KEEN-----STEEQALEDQNAKRAITYQ 434

KK 11

25 Sbjct: 472 DREVEDQDDLDYYESLDKKSKMAKKLRKENHDLERDLIRASRHPELIELGEGDKRGITLD 531

Query: 415 IAKNRGLTPRRKKIDRNPRVKHXXXXXXXXXXXXGQVREVRKEEQR-YSGELSGIRAGVK 473

30 IAKNRGLTPRR K +RNPR+K + + Q Y+GE +GI+AG+ Sbjct: 532

IAKNRGLTPRRPKENRNPRLKKRMRYEKAKKKLASKKAIYKGAPQGGYAGEQTGIKAGLV 591

Score = 80 (12.0 bits), Expect = 6.2e-13, Sum P(2) = 6.2e-13

40 Identities = 29/129 (22%), Positives = 66/129 (51%)

Query: 197 DEAETRVVK-DLAKVSVKEKLKMLRKESP--ELLELIE---DLKVKLTEVKDELEPLLE 249

D ++ + +K D + +++E ++ + + P ELL+++E + ++ L E+

45 ++L+P L

Sbjct: 173 DNSDLKSIKQDSSAAAIEELVQQISPDLPRTELLKILEAKHPEFQLFLDEL-NQLKPQLN 231

Query: 25D LVEQGIIPPGKGSQYLRTKYNLYLNYCSNISFYL-

50 ILKARRVPAHGHPVIERLVTYRNLI 3D8

LV +

Sbjct: 232 EIKEKL-

KTYPSSQLLQAQCTALSTYISFLTFYFALLKDGEEDLKNHPIMVDLVRCKQTW 290

55

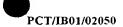
Query: 309 NKLZVVDQKLZ 319 +D+ L+

Sbjct: 291 ESYCGLDEVLT 301



Score = 59 (8.9 bits), Expect = 9.2e-ll, Sum P(2) = 9.2e-llIdentities = 18/59 (30%), Positives = 35/59 (59%) 196 VDEAETRVVKDLAKVSVKEKLKMLRKESPEL---LELIEDLKVKLTEVKDELE--PLLEL 250 ++ DL + E LK+L + PE L+ + LK +L ++E Р E+K++L+ +L 189 IEELVQQISPDLPRT---Sbict: ELLKILEAKHPEFØLFLDELNØLKPØLNEIKEKLKTYPSSØL 245 10 251 VE 252 Query: ++ Sbjct: 246 LQ 247 15 Score = 57 (8.6 bits), Expect = 3.0e-01, Sum P(2) = 2.6e-01Identities = 13/58 (22%), Positives = 26/58 (44%) 367 DSDFDEKAKLKYYKEIEDRQKLKRK--KEENSTEEQALEDQNAKRAITYQIAKNRGLT 422 20 D + +++ L YY+ ++ + K+ +K KE + E RG+T Sbjct: DREVEDADDLDYYESLDKKSKMAKKLRKENHDLERDLIRASRHPELIELGEGDKRGIT 529 25 Score = 42 (6.3 bits), Expect = 5.2e-09, Sum P(2) = 5.2e-09 Identities = 13/51 (25%), Positives = 29/51 (56%) 199 AETRVVKDLAKVSVKEKLKMLRKESPE--LLELIEDLKVKLTEVKDELEPLLE 249 30 + EL +ET + D+++ + FK ++++Z + EL++ + Sbjct: JPD ZELDAIDDIZGMADNZDFKZIKGDZZVVVEFFAGIZEDFA--RTELLKILE 210 35 Score = 39 (5.9 bits), Expect = 1.1e-08, Sum P(2) = 1.1e-08Identities = 8/18 (44%), Positives = 11/18 (61%) 43 YYQDQVDDFHEARSRAAL LO Querv: +Y +Q+D 40 RSRA L 402 FYANQIDQKAAKRSRAVL 419 Sbjct: Pedant information for DKFZphamy2_2f22, frame 3 45 ______ Report for DKFZphamy2_2f22.3 50 **ELENGTHD** 479 54558-00 EMWI 5.50 [[q] TREMBL:SPBC3B8_9 gene: "SPBC3B8.09"; product: EHOMOLI "hypothetical protein": S.pombe chromosome II cosmid c3BB. le-lD 55 YDL153cl le-08 EBLOCKSD PROD528D EBLOCKSI BL003600 Ribosomal protein S9 proteins





	VV () 01/2043	1 01/13/1/2020
5.	EBFOCKZ] EBFOCKZ] EBFOCKZ]	BLOO964A Syndecans proteins PROO624G PROO828H BLOO824B Elongation factor 1 beta/beta//delta chain
	proteins EKW] EKW] EKW]	All_Alpha LOW_COMPLEXITY 24.63 % COILED_COIL 7.10 %
10	SEG ····	SRRRGAAKWAAVRAKAGPTLTDENGDDLGLPPSPGDTSYYQDQVDDFHEARSRAAL •××××××××××××××××××××××××××××××××××
15	COILS	cchhhhhhhhhhhhccccccccccccccccccchhhhhh
	SEG	NEVQSGDEEDGEEEEEEVLALDMDDEDDEDGGNAGEEEEENADDDGGSSVQSEAE ····xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
20	COILS	
25	SEG PRD hccc	PSLSWGQRKKLYYDTDYGSKSRGRQSQQEAEEEEEEEEEEAQIIQRRLAQALQEDDxxxxxxxxxxxxxxxxxxxxxxxxxx
	SEQ FGVA	
30	PRD chhh	hhhhhhhccchhhhhhhhhhhhhhhhhhhhhhhhhhhhh
35		EPLLELVE@GIIPPGKGS@YLRTKYNLYLNYCSNISFYLILKARRVPAHGHPVIER
	PRD hhhhl	hhhhhhhhhhhcccccchhhhhhhhhhhhhhhhhhhhh
40	,	
	SEG ···· PRD hhhhl	RNLINKLSVVD&KLSSEIRHLLTLKDDAVKKELIPKAKSTKPKPKSVSKTSAAACA xxxxxxxxxxxxxxxxxxxxxxxxx
45	COILZ	
	SEQ VTDL	SDDSDFDEKAKLKYYKEIEDRQKLKRKKEENSTEEQALEDQNAKRAITYQIAKNRG
50	COILZ	ccchhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
55	SEG	RKKIDRNPRVKHREKFRRAKIRRRGQVREVRKEEQRYSGELSGIRAGVKKSIKLKxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
	COILZ	

(No Prosite data available for DKFZphamy2_2f22.3)

(No Pfam data available for DKFZphamy2_2f22.3)

DKFZphamy2_2g12

5 group: nucleic acid management

DKFZphamy2_2gl2 encodes a novel 191 amino acid protein with similarity to NVL-2 of Rattus norvegicus.

The novel protein contains 3 EF-hand calcium-binding domains. The related human VILIP Ca-dependend protein specifically binds the 3'-untranslated region of the neurotrophin receptor, trkB, an mRNA localized to hippocampal dendrites in an activity-dependent manner. The new protein exhibists elevated expression in brain and testis.

The new protein can find application in studying the expression profile of brain-specific genes and as a new marker for neuronal cells.

20

strong similarity to NVL-2 (Rattus norvegicus)

Comment for P35332:

25 FUNCTION: MAY BE INVOLVED IN THE CALCIUM-DEPENDENT REGULATION OF RHODOPSIN PHOSPHORYLATION.
TISSUE SPECIFICITY: NEURON-SPECIFIC IN THE CENTRAL AND PERIPHERAL NERVOUS SYSTEM.

MISCELLANEOUS: PROBABLY BINDS TWO OR THREE CALCIUM IONS (BY

30 SIMILARITY)

SIMILARITY: TO OTHER EF-HAND CALCIUM BINDING PROTEINS, BELONGS TO THE RECOVERIN SUBFAMILY.

35 Sequenced by MediGenomix

Locus: /chromosome="l"

Insert length: 4285 bp

40 Poly A stretch at pos. 4258, polyadenylation signal at pos. 4247

1 GGCGGCTCCG GCGCAGACCT TGGAGAGCAC AGCTGCCGGC CCGCGAGCCA
51 GCCTCGGTTC CCGCGGCCCG CCGAGGCTCG GAGCCATCCA GCGACCCGGC
45 101 GACCGGCCTC AGGCCCCGCC ATGGGGAAGA CCAACAGCAA GCTGGCCCCC
151 GAGGTGCTGG AGGACCTTGT TCAGAACACT GAGTTCAGCG AGCAGGAGCT
201 GAAGCAGTGG TACAAGGGCT TCCTGAAGGA CTGCCCCAGC GGCATCCTCA
251 ACCTGGAGGA GTTTCAGCAG CTCTACATCA AGTTCTTCCC CTACGGCGAC
301 GCCTCCAAGT TCGCGCAGCA CGCTTTCCGC ACCTTCGACA AGAACGGCGA
50 351 CGGCACCATC GACTTCCGGG AGTTCATCTG CGCCCTGTCG GTCACCTCCC
401 GCGGCAGCTT CGAGCAGAAG CTCAACTGGG CCTTTGAGAT GTACGACCTG
451 GACGGCGACG GGCGAATCAC GCGCCTGGAG ATGCTGGAGA TCATCGAGGC
501 AATCTACAAG ATGGTGGGCA CCGTGATCAT GATGCGCATG AACCAGGACG
551 GGCTCACGCC CCAGCAGCGT GTGGACAAGA TCTTCAAGAA GATGGACCAG
551 TGACCCATCC ATTGTGTTGC TGCTGCAGTG TGACATGCAG AAGTAGAAGC
701 TGGTGAGGGG CAGGGTCCCT GGCCAGGAG GAGGGACACT CCCAGCCCCCC

```
BD1 TCTCTGGCCC ACCCAGTCCT CTGCCCAAGC CCTTCCTCCC CTCCATCAAG
                             85% ATCTTTGAGG GACCACCTCA CCCTGCAAAA GAGACAGGTC CTCCAGTACC
                             TO TOTAL STANDARDA CONCERN CONCERN CONTROL TO TABLE TO TABLE TO TABLE TO THE TABLE TO TABLE T
                             951 ATAGGGGAGT TGGCTTTTGC CCCAGGAGGT GAGGTTAAGG AGTTGGGGGC
                        LDDL CTGGGGTTCT GGTTAGGAAT TCTCTTGATC CTGGGATTAT GCTTTATAGG
     5
                        1051 ATGTGGTCCC ACAGGCCTGT CACAGGGCCA AATTGGGTCT GTCCATTCCT
                        1101 GAGGCTCCAG ATCCCATAAA GGGGGTCTCT TCCCCATCCC TTCTACTCTA
                        1151 CCTGGCCCTT CCAGCCCCAG CCTTTGGAGC GTTCATTCAG TCCTTTCTTC
                       LIST CCTGGCCCTT CCAGCCCCAG CCTTTGGAGC GTTCATTCAG TCCTTTCTTC
LEDL AGCTAATGAT TACTGAGCAC CTGTTTGGTG CTAAGGATAT GGTCATTTAC
LEDL AAGACACATC TTGTGCCCTC TGGAAGCTCA TAGGGGTTGTG AGGCAAACTT
LEDL CCAGCCGTCA GGGTCTCAGC TAAGCAGAAG GTGCTGGAAG GCTGGTTAGT
LEDL AGGACGAAAT GAAAAGCATT TGGAAGTTTA GGAGCCACGT GAGTGAAAGT
LHSL TTTAAGAAAA ATGAAATTTA TGTCATACTT ATTTTTTAG TACCCTTTAA
 10
                        1501 AGGAGCTACA GTCATTTTAT TATTTCAGGA GGTTAAAATA TACTCTATAT
1551 TACTTGGTTT ATTATAAAAT GATTAAATGA ATAGAGAAAA TATTAATTTT
 15
                     LSDL AGGAGCIACA GICATITIAL TATITAAGGA GGITAAAATA TACTCIATAT

LSDL CAAGGGGAAA AAACCTGAGA AGAAAGGGAG AAAAGACCAT GAAATTTACC

LSDL AGATAACACT TTTTAAGACT AAGTCCTGAG CTGCCACTCT CAGCAGTTTT

L7DL TGCTGCTTCA GCTCTTCTT TTTATTACCT TTTTCAATTC AACAAGCAAC

L7SL TTTCTGCTAC ATACTTACTC CGGTTGGGTG CTGACTTCAG GGACAGGAAA

L8DL AAGCAAGGTT TGCAAAGAGT GAAACTAGTG TATATTCCGT ATCTTGGTAG

L8SL TTCGTTTCTG GATTGGGTTT AGTTTCAGAA CTGGACTTGT TCCTTCACTG

L9DL CCACAGAATC AGAAAGAGCT AGAAGAAAAG GCTCACCTGG CCACTGTTTA

L9SL GGCACCCAGA CATAATTTAT GGACGAAAAG GCTCACCTGG CCACTGTTTA

L9SL GGCCACCAGA CATAATTAT GGACGAAAAG GCTCACCTGG CCACTGTTTA

L9SL GGCCACCAGA CATAATTAT GGACGAAAAG CTTAAAAATG TGCCAGGGAA

L8DL AATTTTGGAT CCAAGGGT ATCAGCCCA AATCTTAGAT CTGCCAGGTA

L9SL TTAATCTTGC TTCTCATCA GGTCTTTCCC CTGTACTTGCT GCCGGAAAAA

LSSL TTAATCTTGC TTCTTCATCA GGTCTTTCTC CTGTACTTGT GATCAGAAAT

L8DL TACCTTTGAC GTGCAGTGAC AGTTGATTTC CTCTTGAACT GCCGGTGAAA

L8DL TCTGGGGGAT GCCAGTGGC CACCTGCT GCTTTCTCC

L9DL TCTGGGGGAT GCCAGTGAC AGTTGATTTC CTCTTGAACT GCCGGTGAAA

L8SL TGCTGACCCC GGGGCAGGGG AATTGCATCT GCTTCTCTC CACGGTCA

L8SL TGGAGGTAGC AAGGCCACTG GGTTGCTCT CCTTTGATCG GGAAAAAC

L8SL TGGAGGTAGC AAGGCCACTG GGTTGCTATC CTCTTTGATCG GGAAAACC

L8SL TGGAGGTAGC AAGGCCACTG GTTGCTATC CTCTTTGATCG GGAAGAACC

L8SL TGGAGGTAGC AAGGCCACTG GTTGCTATC CTCTTTGATCG GGCAGGAACT

L8SL TGGACCACCA TATGGTGAGG CTGGGGAGTT CACATCCTCA GGCAGGAACT

L8SL TGGACCACCA TATGGTGAGG CTGGGGACA TGCCCCACGCCT ACCTGGGCCCT

L8SL TGGCGCACTATT AGGTAGGG CTGCC CACCTCT ACCTGCCCAGGACCATC

L8SL TGCCTCACT TTTTCCCACACCAC TGCCCCACCTCT ACCTGCCCACGACCATCA CACACCATCA CACACCATCA CACACCATCA CACACCATCA CACACCATCA CACACCACCATCA CACACCATCA CACACCATCA CACACCATCA CACACCATCA CACACCATCA CACCACCATCA CACCACCA
 20
 25
 30
 35
                        265 GCTGGCTATT AGGTATGTCT TGTGCGGTCA GTCAGCATCA CAGACACATA
                        2701 GATGCTCACC AGCCTGGCTT AGCTGGGACC TAAATCTTCT GGTGAAAAGC
                       2751 TTTTCACTAA GTGAGGTTCC TTCCCTGCAA ATGCTGAATC TAGCCTAATT
 40
                       2801 CGCAACCACA CAGAATTTCA TGGCTTTCAA AGGCTTGCCA TGTGCCCCAT
                       2851 CTCATTCTAT ACTCACATCC CATGGAGGTG AGGATTTTCA CTTCTTTTCT 29D1 CTAGACTTGG AAGCTGAGAT TCAGAGAGGA AGCATCCCTT GTGCAAGATC
                       2951 ACATAGTCAG GAGGTGACAC AGGGCTAAGA CTTGAACCAA GGCTCTAAGA
                      2951 ACATAGTCAG GAGGTGACAC AGGGCTAAGA CTTGAACCAA GGCTCTAAGA
3001 GGATTTCTTC TTTTCAGAGT CTCTTCCCTG TCCATTTCTG TGACTAAGCT
3051 GTGCAGAGGT TGACAGCAGG GCAAGTTACA TTGATATTCA TCCTTTATAG
3101 GCTTCCTGCT AAAAAGCTTC TGAGATTGTG GTCTTCCAAA AAAAATAGGA
3151 GCTTGGTTGA AGTCCCCACA TTTTCAAGCA CTCAGTGTTC TGCCTCTGGC
3201 AGCTGTGCTA ACAGCTCAGT GCTGTCCTGG GAGTCCTCTG ACTCAGAACC
3251 CTCGAAGCAT CCTGCATTGT CTTTACCCAC CATCATCGTC ACTAAGAGAA
3301 ACATGCCTAC CCATGAAGGC GTGTTTGATT ACTCCAGGCT TCTGGACACA
3351 CATACCCATG GGTGATTTTT GCTCCTCAGG CCCAATATTC TCAGACAGCC
3401 CAGCAGTGTG AACACACAAT GCCAGGCCAG GAACTGGGAC CACCATCTTG
 45
 50
                       3451 CTGATGGAAG GAACAACAGG TGGCCCAGGA CATGCTCCTG CATACTCCTG
                       3501 GGTGTCCCAG GGACTGTGTG CTCAGGAGCA CTGTGGTAGA GCACTGGCCC
55
                       3551 TGCCTTGAGA AGAGACACAG GTCTCCCGTC CCTGCACCAG CTGAGAGAGA
                       3601 CTTGCCACAA AGCACAAGGC TGGCAGAGAT TTATGTATGA CTTGCACAGA
                       3651 CACAAAATA TACAGACAAT CAAAACATTG ATATATTCAA ACTCTCCTTT
```

WO 01/98454

PCT/IB01/02050

3701 AAATTCCAAT CTTATTGCAA CAACTCTGTG AATTGCAAGG TCCCAGAATC
3751 TGCCTTCTCA CATACTCTAC CCTCATTCAT CCTTTTGGGC TAATTGATGA
3801 GCATCTTATT TCTTATCTCT AAAAATTATC AGCAAAGGCT ACTTCAGATG
3851 GCCACTTTAG TCCTTTCAGC TGTAGTCAGG ATTATTTAAC TTACCTGTAT

5 3901 ATCAAAAGTG AAGAAAAAGT TAGTTCATAA GTAAAGGCAC TAAATCCTTT
3951 CCTGACAATG GCAGAGTCTC TAGAGGTAGA AATTTGCCTT GCTGCAGAGA
4001 GAGAAGGAAT GGCGTGGGAT GGGGGAAAGA AAAGAAAGAG AAGAAAGAGA
4051 GAAGCTGGGG TCTCCAGGCA GGGTAGTAAG CTGACACTAA ATATTTTTA
4101 CACAAAAATG TATTGAAGCA ACAAATATTT CCTGAAGATC CACCCTGGGT
10 4151 GAGGCTTTGA GCTGACTTTA GAGAACACTG TGGGGTCAAG AATGTCTTAC
4201 ATGTTTTATT CATCATTCTT GAAAAAAAAA AAAAA

15

BLAST Results

No BLAST result

20

Medline entries

93367470:

25 Kajimoto Ya Shirai Ya Mukai Ha Kuno Ta Tanaka Coa Molecular cloning of two additional members of the neural visinin-like Ca(2+)-binding protein gene family. J Neurochem 1993 Sep;61(3):1091-6

30

35

96079121:

Polymeropoulos M.H., Ide S., Soares M.B., Lennon G.G.; Sequence characterization and genetic mapping of the human VSNLl gene, a homologue of the rat visinin-like peptide RNVPL. Genomics 29(1):273-275(1995).

40

Peptide information for frame 1

ORF from 121 bp to 693 bp; peptide length: 191 Category: strong similarity to known protein 45 Classification: Protein management Prosite motifs: EF_HAND (73-85) EF_HAND (109-121) EF_HAND (159-171)

50

1 MGKTNSKLAP EVLEDLVQNT EFSEQELKQW YKGFLKDCPS GILNLEEFQQ 51 LYIKFFPYGD ASKFAQHAFR TFDKNGDGTI DFREFICALS VTSRGSFEQK 101 LNWAFEMYDL DGDGRITRLE MLEIIEAIYK MVGTVIMMRM NQDGLTPQQR 151 VDKIFKKMDQ DKDDQITLEE FKEAAKSDPS IVLLLQCDMQ K

55

BLASTP hits

PCT/IB01/02050 WO 01/98454

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_2gl2, frame 1 5

No Alert BLASTP hits found

Pedant information for DKFZphamy2_2gl2, frame 1 ______

10

[PIRKW]

Report for DKFZphamy2_2gl2.l

```
ELENGTHD 231
              26277.92
   EMWI
15
    [[q]
              5.26
                    PIR: JHO815 neural visinin-like Ca2+-binding
    [HOMOL]
    protein-type 2 - rat le-107
                                                        ES. cerevisiae 
    EFUNCATD 98 classification not yet clear-cut
   YDR373w3 3e-52
20
    EFUNCATI D3.D1 cell growth ES. cerevisiae, YKL19Dwl 3e-18
    [FUNCAT] [3.07 pheromone response, mating-type determination,
    sex-specific proteins [S. cerevisiae, YKL190w] 3e-18
    EFUNCATD 13.04 homeostasis of other ions ES. cerevisiae.
25
    YKL190w1 3e-18
    EFUNCATD 04.05.01.04 transcriptional control
                                                        ES. cerevisiae,
    YKL190wl 3e-18
             30.03 organization of cytoplasm
                                                   ES. cerevisiae,
    EFUNCATI
    YKLl90wl 3e+l8
   EFUNCATI 11.01 stress response ES. cerevisiae, YGR100w1 7e-04 EBLOCKSI BL00303B S-100/ICaBP type calcium binding protein EBLOCKSI BL00018 EBLOCKSI PR00450G
30
             PR00450F
    EBF0CK23
    EBLOCKSI PROD450E
             PRO0450D
    EBFOCK21
    EBLOCKSI PROD450C
    EBLOCKSI PROD450B
    EBLOCKSI PRODUSOA
                    dlosa__ 1-37-1-5-13 Calmodulin I(Paramecium
40
    E2C0P1
    tetraurelia) 8e-25
                    dlrec__ 1.37.1.5.21 Recoverin Ebovine (Bos
    EZCOPI
    taurus) le-72
    ESCOP3
                   dla4pa_ 1.37.1.2.5 (alcyclin (SLOD) EHuman (Homo
    sapiens), Pl 7e-05
45
                   dlrro___ 1.37.1.4.1 Oncomodulin Erat (Rattus
    [SCOP]
    norvegicus) 2e-17
                   dlsyma_ 1.37.1.2.2 Calcýclin (S100)
                                                           Erat (Rattus
    ESCOPI
    norvegicus) 9e-14
                   d4icb__ 1.37.1.1.1 Calbindin D9K Ebovine (Bos
50
    EZCOPI
    taurus) 2e-18
    ESCOPI
                   dlauib_ 1.37.1.5.19 Calcineurin regulatory subunit
    (B-chain le-45
    EPIRKWI
                  blocked amino end le-99
                   phosphotransferase 3e-D&
55
    [PIRKW]
                   duplication ?e-17
    [PIRKW]
                   tandem repeat 7e-06
    [PIRKW]
                   heterodimer 7e-17
```

```
TPIRKWI
                   heart 7e-0b
                   serine/threonine-specific protein kinase 7e-Ob
    [PIRKW]
    EPIRKUI
                   acetylated amino end 7e-Ob
                   ATP 7e-06
    [PIRKW]
    [PIRKW]
                   skeletal muscle ?e-Ob
                   signal transduction 4e-69
    [PIRKU]
                   protein kinase 3e-08
    IPIRKWD
                   calcium binding le-99
    [PIRKW]
                   alternative splicing le-l3
lipoprotein le-99
    EPIRKWI
    [PIRKW]
10
    [PIRKW]
                   cardiac muscle 7e-06
    EPIRKWI '
                   muscle 7e-0b
                   myristylation le-99
    [PIRKW]
                   EF hand le-99
    [PIRKW]
                   retina le-46
15
    EPIRKWI
    ESUPFAMD calcium-dependent protein kinase 3e-08
ESUPFAMD unassigned calmodulin-related proteins 2e-34
   ESUPFAMI protein kinase homology 3e-DB
ESUPFAMI calmodulin le-99
ESUPFAMI calmodulin repeat homology le-99
20
    EPROSITED EF_HAND
    EPFAMI
                  EF hand
    EKMI
              All_Alpha
    EKWI
              ΒD
25
    SEQ GGSGADLGEHSCRPASQPRFPRPAEARSHPATRRPASGPAMGKTNSKLAPEVLEDLVQNT
    lrec-
         30
    SEQ EFSEQELKQWYKGFLKDCPSGILNLEEFQQLYIKFFPYGDASKFAQHAFRTFDKNGDGTI
             СССИНИНИНИНИНИНИТТТТЕЕЕНИНИНИНИНТТТТСИНИНИНИНИНЫ—--
    lrec−
    --CEE
35
    SEQ DFREFICALSVTSRGSFEQKLNWAFEMYDLDGDGRITRLEMLEIIEAIYKMVGTVIMMRM
    lrec-
         SEQ NQDGLTPQQRVDKIFKKMDQDKDDQITLEEFKEAAKSDPSIVLLLQCDMQK
              ТТТТТСНННННННННННССТТТТЕЕСННННННННННСНННННННСССННН
40
    lrec-
                         Prosite for DKFZphamy2_2gl2.1
45
                             EF_HAND
                                                       PDOCDDDAB
    810002A
                  JJ3->J5P
                             EF_HAND
    810002A
                  149->162
                                                       PDOCUUDA
                             EF_HAND
    810002A
                  199->212
                                                       PDOCOOOT9
50
                           Pfam for DKFZphamy2_2q12.1
    HMM_NAME EF hand
55
    HMM
                       *EIgEMFrmMDkDGDGyIDFEEFmeMMkem*
                          Q +FR +DK+GDG+IDF EF+ +++
```

WO 01/98454 PCT/IB01/02050 104 FAQHAFRTFDKNGDGTIDFREFICALSVT 735 Query 27.15 140 7P8 ŀ 29 dkfzphamy2_2gl2.1 strong similarity to NVL-2 (Rattus norvegicus) 5 Alignment to HMM consensus: *ElqEMFrmMDkDGDGyIDFEEFmeMMkem* Query ++++F+M+D DGDG+I+ E++E++ ++ 140 KLNWAFEMYDLDGDGRITRLEMLEIIEAI dkfzphamy2 768 l 29 dkfzphamy2_2gl2.l strong 10 Query 278 similarity to NVL-2 (Rattus norvegicus) Alignment to HMM consensus: *EIgEMFrmMDkDGDGyIDFEEFmeMMkem* ++++F++MD+D+D +I+ EEF+E+ K+ 190 RVDKIFKKMDQDKDDQITLEEFKEAAKSD 15 Query 238 .

5 group: amygdala derived

DKFZphamy2_2il7 encodes a novel 462 amino acid protein without similarity to known proteins.

- 10 Most ESTs are derived from brain and pancreas. No informative BLAST results; No predictive prosite, pfam or SCOP motife.
- The new protein can find application in studying the expression profile of amygdala-specific genes.

unknown protein

20 perhaps complete cds.

Sequenced by MediGenomix

Locus: unknown

25

Insert length: 3473 bp
Poly A stretch at pos. 3454, polyadenylation signal at pos. 3436

30 L GATATCCCAA TCTTTGGACT GCATCCTGGT TGCCTCTACT GTGGTCACCT 51 TTGGGAAGAA ATGTCTTCTG TAAAAAGAAG TCTGAAGCAA GAAATAGTTA LOL CTCAGTTTCA CTGTTCAGCT GCTGAAGGAG ATATTGCCAA GTTAACAGGA 151 ATACTCAGTC ATTCTCCATC TCTTCTCAAT GAAACTTCTG AAAATGGCTG 201 GACTGCTTTA ATGTATGCGG CAAGGAATGG GCACCCAGAG ATAGTCCAAT 251 TTCTGCTTGA GAAAGGGTGT GACAGATCAA TTGTCAATAA ATCAAGGCAG 35 301 ACTGCACTGG ATATTGCTGT ATTTTGGGGT TATAAGCATA TAGCTAATTT 351 ACTAGCTACT GCTAAAGGTG GGAAGAAGCC TTGGTTCCTA ACGAATGAAG 401 TGGAAGAATG TGAAAATTAT TTTAGCAAAA CACTACTGGA CCGGAAAAGT 451 GAAAAGAGGA ATAATTCTGA CTGGCTGCTA GCTAAAGAAA GCCATCCAGC 40 501 CACAGTTTTT ATTCTTTTCT CAGATTTAAA TCCCTTGGTT ACTCTAGGTG 551 GCAATAAAGA AAGTTTCCAA CAGCCAGAAG TTAGGCTTTG TCAGCTGAAC LOD TACACAGATA TAAAGGATTA TTTGGCCCAG CCTGAGAAGA TCACCTTGAT 651 TTTTCTTGGA GTAGAACTTG AAATAAAAGA CAAACTACTT AATTATGCTG 701 GTGAAGTCCC GAGAGAGGAG GAAGATGGAT TGGTTGCCTG GTTTGCTCTA 45 751 GGTATAGATC CTATTGCTGC TGAAGAATTC AAGCAAAGAC ATGAAAATTG BOD TTACTTTCTT CATCCTCCTA TGCCAGCCCT TCTGCAATTG AAAGAAAAAG **B51 AAGCTGGGGT TGTAGCTCAA GCAAGATCTG TTCTTGCCTG GCACAGTCGA** 901 TACAAGTTTT GCCCAACCTG TGGAAATGCA ACTAAAATTG AAGAAGGTGG 951 CTATAAGAGA CTATGTTTAA AAGAAGACTG TCCTAGTCTC AATGGCGTCC 951 CTATAAGAGA CTATGTTTAA AAGAAGACTG TCCTAGTCTC AATGGCGTCC
1001 ATAATACCTC ATACCCAAGA GTTGATCCAG TAGTAATCAT GCAAGTTATT
1051 CATCCAGATG GGACCAAATG CCTTTTAGGC AGGCAGAAAA GATTTCCCCC
1101 AGGCATGTTT ACTTGCCTTG CTGGATTTAT TGAGCCTGGA GAGACAATAG
1151 AAGATGCTGT TAGGAGAGAA GTAGAAGAGG AAAGTGGAGT CAAAGTTGGC
1201 CATGTTCAGT ATGTTGCTTG TCAACCATGG CCAATGCCTT CCTCCTTAAT
1251 GATTGGTTGC TTAGCTCTAG CAGTGTCTAC AGAAATTAAA GTTGACAAGA
1301 ATGAAATAGA GGATGCCCGC TGGTTCACTA GAGAACAGGT CCTGGATGTT
1351 CTGACCAAAG GGAAGCAGCA GGCATTCTTT GTGCCACCAA GCCGAGCTAT
1401 TGCACATCAA TTAATCAAAC ACTGGATTAG AATAAATCCT AATCTCTAAA 50 55

	1451	TCTAAGAACT	AAGCTTTGAG	TATTATTAA	TAATTTCTAA	TAACACTCAT
	1501	TCCTCAAGTG	ATATTAGAGA	TTATTCAGTA	CTCTTGAGAG	TGTCACAACA
	1551	CAAAATACGA	TGTTGGGTTT	TCGAAATATT	TTCAAAGTGT	TCTGTCTTAA
	1601	TCACAAATTC	ATATTTTTAC	ACATTTTTAC	AATATTGCCT	CAGATTATGT
5	1651	TAAATTTGGG	TCAGTCTTCT	CTGAACTTTT	TCTCTCTCGG	TTTCTTTTCT
_	1701	TCCTTCACAG	TTTTATCTCA	CAAAACCATT	TTTCTAATAA	GAGACATCAT
	1751	GTTGGAAAGA	TGTTGTAGAA	ATGTGCATAA	ATTTCAGTGC	CTCTTGTAAG
	1801	CATTAAACTG	ATGATGAAGA	AAGTTCCTGA	TTTGAGAAAT	GAATCAAAGT
	1851	AATTTTAATG	AATTTTTAGC	TTGTATTAGC	TTGAGTTAGC	TGGCATTGAT
10	1907	TTTTTAGTCC	TTTTGTTACC	TTTAAGTTGT	CAATATATGG	TTTTTGTTCA
10	1951	TCTCCCCATT	GTAGTCCCAC	TTGCTCTTTC	CTGGGGGTTC	CATTGTTCTA
	2007	GCAGTGGAGG	TGTTACAGTG	TCGCCACTCG	TCTAATTTGA	CCAGTGTTAA
	2051	GAATTTTCTA	ATTTAATAAT	TTAATAGTGA	TCTCAATACC	ACACCCTCAT
	5707	GGAAGGAGAA	AAGCATACTA	TTATATCTGG	GACCTCTCTT	TTAGACCTAA
15	2727	AATTAATTAA	CATATCTACT	TATATGTTAC	TTATACCTAA	AGCTGTTATT
13	5507	AAGACAAACC	AAGATTCTCT	GCTTTTGCAC	TGAAATTAAA	CTTGAAAGGA
	2251	ATTCTCCTCA	AAGGTCGGAT	ATTAAATAAG	TCCCAGGCAG	ATTTACATAT
	5307	TTAATTTAAA	ACATTGGCTT	TATTTCATTT	TGTGATGAGT	GATGTATCTG
	2351	TGTTAACAAA	AAATTGTATA	ATCATTACCA	ATACTATTTA	TTATGCTCAA
20	2401	ATATATCTTG	GCTTTGACCT	TATTTCAACA	CATTCTAAGA	AGCCTTGACA
20	2451	AAGTAAGTAT	ATTTTAGAGC	TGAATCAGTA	AGATTCTAGA	GAAAGCAAAA
	2501	CATAGTAGTT	CACAATTTTG	CAACATAGAA	AGTCACATTT	TGAAAGGCTA
	2551	TTTTGAAATT	GATTTAATAG	CTATTATAGT	TTATGAATAT	CAAAATTTGT
	5227	ATAATTTGCA	TCTTTACTAA	TGTATGCTAG	AGCTACAAGA	GACCTTAAGG
25	5P27	ATAATATATG	AAATTAGCTT	TCCTTATTTT	ATAGATAAGG	AAAAGAAAT
23	5207	TGTGAAAGGT	GAATTTACCT	AATTAGTGAA	AGTTACATAA	CTAATTACAA
	2751	CAGTCTGTAC	TATATAATGC	AGAGGACGAT	TCTCCCTGTA	AAAGGAACTA
	5907	GAAGCTATTA	CTAAAAATAT	ATATAGACAA	AATTAAAAGA	AGGAATGATA
	2851	AGAATAAATT	TAATTTACCA	AATATTGTTA	ATTAAAATTT	TAGATACTTA
30	5407	ACATTTATTT	AACTTAAATA	AAAGATAACT	GTCAGATAAA	ACTTTATTTT
50	2951	ACTAATGAGC	AGTGATTTTC	TTAGGAATTG	ATGAAGGCTT	ATTGGTATCA
	3007	AGAATTTAAA	CCAAATTAAA	ACTGACAGAG	GACATTTAGA	TACATAATAA
	3051	AATTCGAGCT	ACATAAGTAT	ATGGAAAATA	ATGTACCTTG	ATTATTATGA
	3707	AATAGAGCAT	CTTGAAATTC	AGTTTTACTC	TAAATGTACT	TTTAATACTT
35	3151	GCAGATTCTA	AGATTACATT	GTGAAATTCC	AGGTTTTCAT	AATGTTAAAA
55	3507	TAGGAAAGTA	GAATATAAAG	TATCAACAAG	TGTAGTTATA	CATTTTGTTT
	3251	TGGATATTTA	ATCCTTACTT	GGGAAAAAT	CAGCATCTAG	GTAAATTATT
	3307	ATTTTAATAA	GAACTCTTAA	ATTGCCAACC	TCTGAGAGGT	GAAAAGCTAT
	3351	GTAAATAGAA	GGAATGGCCA	GTTCAAAAGA	ATAGTAGAAG	TGATAGTGCC
40	3401	GTGAATGTAT	TCTACTGGAA	ATGAATGTAA	TAATACATTA	AATTTTTAAA
70	3451	ATCGAAAAA	AAAAAAAAA	AAA		
	ענירי	A I COUNTANT				

BLAST Results

45

No BLAST result

50

Medline entries

No Medline entry

55

Peptide information for frame l

ORF from 61 bp to 1446 bp; peptide length: 462 Category: putative protein Classification: unclassified

5 Prosite motifs: MUTT (355-374)

1 MSSVKRSLKØ EIVTØFHCSA AEGDIAKLTG ILSHSPSLLN ETSENGUTAL
51 MYAARNGHPE IVØFLLEKGC DRSIVNKSRØ TALDIAVFUG YKHIANLLAT
10 101 AKGGKKPUFL TNEVEECENY FSKTLLDRKS EKRNNSDULL AKESHPATVF
151 ILFSDLNPLV TLGGNKESFØ QPEVRLCØLN YTDIKDYLAØ PEKITLIFLØ
172 PELEIKDKLL NYAGEVPREE EDGLVAUFAL GIDPIAAEEF KØRHENCYFL
173 HPPMPALLØL KEKEAGVVAØ ARSVLAUHSR YKFCPTCGNA TKIEEGGYKR
174 BOL LCLKEDCPSL NGVHNTSYPR VDPVVIMØVI HPDGTKCLLG RØKRFPPGMF
175 BSL TCLAGFIEPG ETIEDAVRRE VEEESGVKVG HVØYVACØPU PMPSSLMIGC
176 HOL LALAVSTEIK VDKNEIEDAR UFTREØVLDV LTKGKØØAFF VPPSRAIAHØ
177 HSL STANDARDEN VANDENDER VANDENDER VEEESGVKVG HVØYVACØPU PMPSSLMIGC
175 HOL LALAVSTEIK VDKNEIEDAR UFTREØVLDV LTKGKØØAFF VPPSRAIAHØ
175 LIKHUIRINP NL

20

BLASTP hits

No BLASTP hits available

25 Alert BLASTP hits for DKFZphamy2_2il7, frame l

No Alert BLASTP hits found

Pedant information for DKFZphamy2_2il7, frame 1

Report for DKFZphamy2_2il7.1

35 .ELENGTHD 462

EMWD 52076-25

EpID 6.38

CHOMOLD TREMBL:SPBC1778_3 gene: "SPBC1778.03c"; product:
"conserved hypothetical protein"; S.pombe chromosome II cosmid

40 cl778. le-45

45 [FUNCAT] I genome replication, transcription, recombination and repair [M. jannaschii, MJ]]49 nucleotide pyrophosphohydrolase] le-04

EBLOCKSI BLOOZI9F Anion exchangers family proteins

EBFOCKZI BF07543B

50 EBLOCKSI DMD1909

EBLOCKZI PF000534

IBLOCKSI BLOOM93 mutT domain proteins

ISCOPI dlawcb_ 1.91.3.1.2 GA binding protein (GABP) alpha GA bindini 2e-35

55 ESUPFAMI hypothetical protein HIO432 le-22 EPROSITEI MUTT L

EPFAMD Bacterial mutT protein

[PFAM] Ank repeat

IKWD Irregular
IKWD 3D

SEQ MSSVKRSLKQEIVTQFHCSAAEGDIAKLTGILSHSPSLLNETSENGWTALMYAARNGHPE TTEETTTEEHHHHHHHHHCCHH SEQ IVQFLLEKGCDRSIVNKSRQTALDIAVFWGYKHIANLLATAKGGKKPWFLTNEVEECENY 10 lawcB HHHHHHHHCCTTTTCBTTTBCHHHHHHHHHCCHHHHHHH......... SEQ FSKTLLDRKSEKRNNSDWLLAKESHPATVFILFSDLNPLVTLGGNKESFQQPEVRLCQLN lawcB . 15 SEQ YTDIKDYLAQPEKITLIFLGVELEIKDKLLNYAGEVPREEEDGLVAWFALGIDPIAAEEF lawcB 20 SEQ K@RHENCYFLHPPMPALL@LKEKEAGVVA@ARSVLAWHSRYKFCPTCGNATKIEEGGYKR lawcB SEQ LCLKEDCPSLNGVHNTSYPRVDPVVIMQVIHPDGTKCLLGRQKRFPPGMFTCLAGFIEPG 25 lawcB SEQ ETIEDAVRREVEEESGVKVGHVQYVACQPWPMPSSLMIGCLALAVSTEIKVDKNEIEDAR 30 lawcB SEQ WFTREQVLDVLTKGKQQAFFVPPSRAIAHQLIKHWIRINPNL lawcB 35 Prosite for DKFZphamy2_2il7.1 40 EPBOOZ9 355->375 MUTT PD0C00695 Pfam for DKFZphamy2_2il7.1 45 HMM_NAME Ank repeat HMM *GyTPLHIAARyNNvEMVrlLLQHGADIN* 50 G+T+L++AAR+++ E+V++LL++G D 46 GWTALMYAARNGHPEIVQFLLEKGCDRS Query 73 HMM_NAME Bacterial mutT protein 55 *ILMiqRedppnHYdtHhgdWIFPGGkIEeGETPEQCarREIWEETGI*

L+++++ +++
++G+IE+GET+E+++RRE++EE+G+
Query 337 CLLGRQKRF--PPG---MFTCLAGFIEPGETIEDAVRREVEEESGV 377

DKFZphamy2_2ol3

5 group: intracellular transport and trafficing

DKFZphamy2_2013 encodes a novel 590 amino acid protein with high similarity to murine synaptotagmin 3.

- 10 The novel protein contains two C2 domains. The C2 domain is thought to be involved in calcium-dependent phospholipid binding. Synaptotagmins are essential for Ca(2+)-regulated exocytosis of neurosecretory vesicles.
- 15 The new protein can find application in modulating/blocking synaptic activity.

similarity to synaptotagmin 3 (Mus musculus)

20 Sequenced by MediGenomix

Locus: unknown

25 Insert length: 2931 bp
Poly A stretch at pos. 2912, polyadenylation signal at pos. 2884

1501 AGGCTCTGGA GAGGCAGGCA CAGGGGCACC CTGTGGCCGT ATCAGCTTCG 1551 CCCTGCGGTA CCTCTATGGC TCGGACCAGC TGGTGGTGAG GATCCTGCAG JUDI GCCCTGGACC TCCCTGCCAA GGACTCCAAC GGCTTCTCAG ACCCCTACGT 1651 CAAGATCTAC CTGCTGCCTG ACCGCAAGAA AAAGTTTCAG ACCAAGGTGC 5 1701 ACAGGAAGAC CCTGAACCCC GTCTTCAATG AGACGTTTCA ATTCTCGGTG 1751 CCCCTGGCCG AGCTGGCCCA ACGCAAACTG CACTTCAGCG TCTATGACTT 1801 TGACCGCTTC TCGCGGCACG ACCTCATCGG CCAGGTGGTG CTGGACAACC 1851 TCCTGGAGCT GGCCGAGCAG CCCCCTGACC GCCCGCTCTG GAGGGACATC 1901 GTGGAGGGCG GCTCGGAAAA AGCAGATCTT GGGGAGCTCA ACTTCTCACT 1951 CTGCTACCTC CCCACGGCCG GGCGCCTCAC CGTGACCATC ATCAAAGCCT 10 2001 CTAACCTCAA AGCGATGGAC CTCACTGGCT TCTCAGACCC CTACGTGAAG 2D53 GCCTCCCTGA TCAGCGAGGG GCGGCGTCTG AAGAAGCGGA AAACCTCCAT 2101 CAAGAAGAAC ACGCTGAACC CCACCTATAA TGAGGCGCTG GTGTTCGACG 2151 TGGCCCCGA GAGCGTGGAG AACGTGGGGC TCAGCATCGC CGTGGTGGAC 2201 TACGACTGCA TCGGGCACAA CGAGGTGATC GGCGTGTGCC GTGTGGGCCC 15 2251 CGACGCTGCC GACCCGCACG GCCGCGAGCA CTGGGCAGAG ATGCTGGCCA Z3DJ ATCCCCGCAA GCCCGTGGAG CACTGGCATC AGCTAGTGGA GGAAAAGACT 2351 GTGACCAGCT TCACAAAAGG CAGCAAAGGA CTATCAGAGA AAGAGAACTC 2401 CGAGTGAGGG GTCTGGCCTA GGCCCGGGAT CGGACCAGGC TCCCTCAGGA 2451 CCCCATCCTT TCCTGCCCGG ACCGTGAATT CATCTCCTTG AAGCCATAAC 20 2503 GTCCGAGCTG CTGGTGCGGG GCAGCCCTGG CCCTAGGCTT CCTAACCCTG 2553 GAAGCGAGAG GATGAGAGGA GGCCGGCCCA GCTCCTTCTT TCAGGGTGGG 2603 GGTCATTCAG CCTCCACTGT GTCTGTCTTT TCTTCCCTGG GGCTCCCCCT 2653 CGAGGCGAGG GGCCATGCAT GTCTGGGGGA CCCCTGCCCC CCAAAACCCT 25 2703 CTGTCTGTCT CTGTCTCTTT GCTGTTTGTC CAAGACTCAG TGTCCCGACC 2751 CTTGTTCTCG CCGTGAATGT CAATGGGCCA ATCCTCTCTG TCCTTTCAGA 2BO3 CACACACA CCTGTGTCCA CCCCTTCTGT TCGCCACACC CTGCGTCTGG 2851 CCGGTCCCCC CACTGCTGCT GCTATCAACG CCAGAATAAA CACACTCTGT 2901 GGGTCTCACT CCAAAAAAA AAAAAAAAA A

30

BLAST Results

35 Entry MMABA93_1 from database TREMBL:
product: "synaptotagmin 3"; Mus musculus mRNA for synaptotagmin 3;
complete cds.

Score = 1814, P = 5.7e-239, identities = 362/450, positives = 369/450, frame +2

45

40

Medline entries

96064733:

Fukuda M. Kojima T. Aruga J. Niinobe M. Mikoshiba K.: Functional diversity of C2 domains of synaptotagmin family. Mutational analysis of inositol high polyphosphate binding domain. J. Biol Chem 1995 Nov 3:270(44):26523-7

ORF from 635 bp to 2404 bp; peptide length: 590 Category: strong similarity to known protein 5 Classification: Cell signaling/communication Prosite motifs: C2_DOMAIN_1 (323-338) C2_DOMAIN_1 (455-470)

10 J MSGDYEDDLC RRALILVSDL CARVRDADTN DRCQEFNDRI RGYPRGPDAD
51 ISVSLLSVIV TFCGIVLLGV SLFVSWKLCW VPWRDKGGSA VGGGPLRKDL
101 GPGVGLAGLV GGGGHHLAAG LGGHPLLGGP HHHAHAAHHP PFAELLEPGS
151 LGGSDTPEPS YLDMDSYPEA AAAAVAAGVK PSQTSPELPS EGGAGSGLLL
201 LPPSGGGLPS AQSHQQVTSL APTTRYPALP RPLTQQTLTS QPDPSSEERP
15 251 PALPLPLPGG EEKAKLIGQI KPELYQGTGP GGRRSGGGPG SGEAGTGAPC
301 GRISFALRYL YGSDQLVVRI LQALDLPAKD SNGFSDPYVK IYLLPDRKKK
351 FQTKVHRKTL NPVFNETFQF SVPLAELAQR KLHFSVYDFD RFSRHDLIGQ
401 VVLDNLLELA EQPPDRPLWR DIVEGGSEKA DLGELNFSLC YLPTAGRLTV
451 TIIKASNLKA MDLTGFSDPY VKASLISEGR RLKKRKTSIK KNTLNPTYNE
20 501 ALVFDVAPES VENVGLSIAV VDYDCIGHNE VIGVCRVGPD AADPHGREHW
551 AEMLANPRKP VEHWHQLVEE KTVTSFTKGS KGLSEKENSE

25 BLASTP hits

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_2ol3, frame 2

TREMBL:MMABA93_1 product: "synaptotagmin 3"; Mus musculus mRNA for synaptotagmin 3; complete cds.; N = 2; Score = 1814; P = 1.1e-239

35

>TREMBL:MMAB893_1 product: "synaptotagmin 3"; Mus musculus mRNA for

synaptotagmin 3, complete cds. Length = 557

40 HSPs:

55

Score = 1814 (272.2 bits), Expect = 1.1e-239, Sum P(2) = 1.1e-239

45 Identities = 362/449 (80%), Positives = 369/449 (82%)

FAELLEPG LGGS+ PEPSYLDMDSYPE GVKPSQT

50 Sbjct: 143
FAELLEPGGLGGSELPEPSYLDMDSYPEAAVASVVAAGVKPSQTSPELPSEGGTGSGLLL 202

Query: 201
XXXXXXXXXXXQSHQQVTSLAPTTRYPALPRPLTQQTLTSQPDXXXXXXXXXXXXXXXXXX 260

QSHQQVTSLAPTTRYPALPRPLTQQTLT+Q D
Sbjct: 203
LPPSGGGLPSAQSHQQVTSLAPTTRYPALPRPLTQQTLTTQADPSTEERPPALPLPLPGG 262

Query: 261

PCGRISFALRYLYGSDQLVVRI

5 Sbjct: 263 EEKAKLIGQIKPELYQGTGPGGRRGGGSGEAGA----PCGRISFALRYLYGSDQLVVRI 317

Query: 321

LQALDLPAKDSNGFSDPYVKIYLLPDRKKKFQTKVHRKTLNPVFNETFQFSVPLAELAQR 380

10

LQALDLPAKDSNGFSDPYVKIYLLPDRKKKFQTKVHRKTLNP+FNETFQFSVPLAELAQR
Sbjct: 318
LQALDLPAKDSNGFSDPYVKIYLLPDRKKKFQTKVHRKTLNPIFNETFQFSVPLAELAQR 377

15 Query: 381

KLHFSVYDFDRFSRHDLIGQVVLDNLLELAEQPPDRPLWRDIVEGGSEKADLGELNFSLC 440

KLHFSVYDFDRFSRHDLIGQVVLDNLLELAEQPPDRPLURDI+EGGSEKADLGELNFSLC Sbjct: 378

20 KLHFSVYDFDRFSRHDLIGQVVLDNLLELAEQPPDRPLWRDILEGGSEKADLGELNFSLC 437

Query: 441
YLPTAGRLTVTIIKASNLKAMDLTGFSDPYVKASLISEGRRLKKRKTSIKKNTLNPTYNE 500

25 YLPTAGRLTVTIIKASNLKAMDLTGFSDPYVKASLISEGRRLKKRKTSIKKNTLNPTYNE
Sbjct: 438
YLPTAGRLTVTIIKASNLKAMDLTGFSDPYVKASLISEGRRLKKRKTSIKKNTLNPTYNE 497

Query: 501

30 ALVFDVAPESVENVGLSIAVVDYDCIGHNEVIGVCRVGPDAADPHGREHWAEMLANPRKP 560

ALVFDVAPESVENVGLSIAVVDYDCIGHNEVIGVCRVGP+AADPHGREHWAEMLANPRKP Sbjct: 498 ALVFDVAPESVENVGLSIAVVDYDCIGHNEVIGVCRVGPEAADPHGREHWAEMLANPRKP 557

35

Query: 561 VEHWHQLVEEKTVTSFTKGSKGLSEKENSE 590 VEHWHQLVEEKT++SFTKG KGLSEKENSE 587 Spjct: 558 VEHWHQLVEEKTLSSFTKGGKGLSEKENSE 587

40 Score = 520 (78.0 bits), Expect = 1.le-239, Sum P(2) = 1.le-239
Identities = 98/100 (98%), Positives = 99/100 (99%)

Query: L MSGDYEDDLCRRALILVSDLCARVRDADTNDRCQEFND-RIRGYPRGPDADISVSLLSVI 59

45 MSGDYEDDLCRRALILVSDLCARVRDADTNDRCQEFN+ RIRGYPRGPDADISVSLLSVI

Sbict: 1

MSGDYEDDLCRRALILVSDLCARVRDADTNDRCQEFNELRIRGYPRGPDADISVSLLSVI 60

50 Query: 6D VTFCGIVLLGVSLFVSWKLCWVPWRDKGGSAVGGGPLRKD 99 VTFCGIVLLGVSLFVSWKLCWVPWRDKGGSAVGGGPLRKD

Sbjct: bl VTFCGIVLLGVSLFVSWKLCWVPWRDKGGSAVGGGPLRKD 100

55 Pedant information for DKFZphamy2_2ol3, frame 2

Report for DKFZphamy2_2013.2

```
590
    ELENGTHI
              63304-02
    EMMD
5
              6.16
    [[q]
                   TREMBL:MMABA93_1 product: "synaptotagmin 3"; Mus
    [HOMOL]
    musculus mRNA for synaptotagmin 3, complete cds. D.O
              99 unclassified proteins
    EFUNCATI
                                            IS. cerevisiae, YML072cl
    be-10
    EFUNCATI
              Ol.Ob.Ol lipid, fatty-acid and sterol biosynthesis
10
       ES. cerevisiae, YGR170w1
    EFUNCATI
              30.08 organization of golgi
    7e-06
              BLO1224A N-acetyl-gamma-glutamyl-phosphate reductase
    EBFOCK2
15
    proteins
              BLOLOIB 0xysterol-binding protein family proteins
    [BLOCK2]
    EBFOCKZ
              PF01368B
    [SCOP]
                   dla25a_ 2.6.1.2.2 C2 domain from protein kinase c
    (beta)
            ERa 2e-27
20
    EGCOPI
                   dlrsy___ 2.6.1.2.1 Synaptogamin I, first C2 domain
    ERat (Rattu 4e-43
                   dlrlw___ 2.6.1.1.2 A domain from cytosolic
    phospholipase A2 [Huma 5e-12
    [COD]
                   dlgasb2 2.6.1.1.1 Phosphoinositide-specific
25
    phospholipase C 4e-27
    CPIRKW]
                   phosphotransferase 7e-15
    [PIRKW]
                   duplication be-76
    EPIRKU
                   synaptic vesicle le-167
                   phorbol ester binding 2e-14
    [PIRKW]
                   zinc 2e-14
30
    EPIRKUJ
                   transmembrane protein [].[]
    [PIRKW]
    [PIRKW]
                   serine/threonine-specific protein kinase 7e-15
    EPIRKW
                   membrane trafficking 0.0
                   phospholipid binding be-76
    EPIRKWI
35
                   autophosphorylation 7e-15
    [PIRKW]
                   ATP 7e-15
    [PIRKW]
    [PIRKW]
                   phosphoprotein 7e-15
    EPIRKUI
                   glycoprotein le-167
                   calcium binding 5e-34
    EPIRKUJ
40
                   alternative splicing Le-10
    EPIRKWI
                   dimer le-75
    EPIRKUJ
    EPIRKUJ
                   membrane protein le-167
                   calmodulin binding 2e-74
    EPIRKWI
    ESUPFAMD ras-specific GAP catalytic domain homology le-OB
             protein kinase C zinc-binding repeat homology 7e-15
45
    ESUPFAMI
             protein kinase homology 7e-15
    CSUPFAMD
              protein kinase C alpha 7e-15
    ESUPFAM3
             HsC2 phosphatidylinositol 3-kinase le-09
    ESUPFAMD
              synaptotagmin 0.0
    ESUPFAM3
             PX domain homology le-09
50
    ESUPFAM3
              pleckstrin repeat homology le-OB
    ESUPFAMD
              protein kinase C C2 region homology D.O
    ESUPFAMD
    EPROSITED C2_DOMAIN_L
    EPFAMJ
                   C2 domain
55
    [KW]
              Irregular
    EKWJ
             LOW_COMPLEXITY
                                20.00 %
```

	Jrsy- SEG SEQ						RGPDADISVSLLSVIV
5	JI Sy-	• • • • • • • • • • • • • • • • • • • •	• • • • • • • • • • • • • • • • • • • •				• • • • • • • • • • • • • • • • • • • •
	SEQ SEG lrsy-	-				×××××	SLAGLVGGGGHHLAAG «×××××××××××××
10		• • • • • • • •	• • • • • • • • • •	• • • • • • • • •	••••••	• • • • •	••••••
15	zeg zeg zeg	_xxxxxxx -		<×ו••••	• • • • • • • • • •		OSYPEAAAAAVAAGVK
15						• • • • •	•••••••
20	SEQ SEG Irsy-	×××× -	*****	(XXXXXXXX	xxx	• • • • •	RYPALPRPLTQQTLTS
20							
25	SEQ Lrsy-	xxxx -	*****	«××····	· · · · · · · · × ×	(xxxxx)	GGGPGSGEAGTGAPC
25							••••••
	SEQ SEG lrsy-						PDRKKKFQTKVHRKTL
30		CEEEEEEE	EETTTTEEEEE	EEEEECCCC	CBTTTBBCEEE	EEEEET	TTTTTEECCCTTTBT
	SEQ SEG lrsy-						ILLELAE@PPDRPLWR
35	,	TTEEEEEE	ЕЕСССННННН	CEEEEEEE	CTTTTCCEEEE	Έ••••	• • • • • • • • • • • • • • • • • • • •
40	SEQ SEG lrsy-				TVTIIKASNLK		FSDPYVKASLISEGR
40						• • • • • •	
45	SEQ SEG lrsy-	• • • • • • •			• • • • • • • • •		IGHNEVIGVCRVGPD
70	55.4						
	SEQ AADPHGREHWAEMLANPRKPVEHWHQLVEEKTVTSFTKGSKGLSEKENSE SEG						
50							
			Pro	site for	DKFZphamy2_	2013.ē	!
55	40029 40029		323->339 455->471	IAMOD_SO			PD0000380 PD0000380

Pfam for DKFZphamy2_2ol3.2

5 HMM_NAME C2 domain *LtVrIIeARNLWkMDMnGfSDPYVKVdMdPdpkDtkKWKTkTiWNNGLN L+VRI++A +L+++D+NGFSDPYVK++++PD+K 10 KK++TK++++ LN Query 316 LVVRILQALDLPAKDSNGFSDPYVKIYLLPDRK--KKFQTKVHRKT-LN 361 HMM PVWNEEeFvFedIPyPdlqrkMLRFaVWDWDRFSRBDFIGHCi* 15 PV+N E+F+F +P+ +L+ + L+F+V+D+DRFSR+D+IG+++ 362 PVFN-ETFQFS-VPLAELAQRKLHFSVYDFDRFSRHDLIGQVV Query 402 20 *LtVrIIeARNLWkMDMnGfSDPYVKVdMdPdpkDtkKWKTkTiWNNGLN LTV+II+A NL++MD +GFSDPYVK +++ + +++KK+KT+++N+ LN Query 448 -LTVTIIKASNLKAMDLTGFSDPYVKASLISEGRRLKKRKTSIKKNT-LN 495 25 HMM PVWNEEeFvFedIPyPdlgrkMLRFaVWDWDRFSRBDFIGHCi* P++N E +VF+ ++ ++ +++ L +AV D+D++++++1G+C+ 496 PTYN-EALVFD-VAPESVENVGLSIAVVDYDCIGHNEVIGVCR Query 536 30

WO 01/98454

5

20

25

PCT/JB01/02050

DKFZphamy2_7j5

group: differentiation/development

DKFZphamy2_7j5 encodes a novel 693 amino acid protein with similarity to Tspyll testis-specific Y-encoded-like protein of Mus musculus.

TSPY genes are arranged in clusters on the Y chromosome of many mammalian species. TSPY is believed to function in early spermatogenesis and is a candidate for GBY, the putative gonadoblastoma-inducing gene on the Y. The TSPY family forms part of a superfamily, TTSN, with autosomal representatives, highly conserved in mammals and beyond.

The new protein can find application in studying the expression profile of testis- and brain-specific genes and diagnosis/therapy of malfunctioning male fertility.

HRIHFB2216

similarity to Y-linked Gene of Mus musculus

Sequenced by BMFZ

Locus: unknown

30 Insert length: 2819 bp
Poly A stretch at pos. 2800, polyadenylation signal at pos. 2779

L AGGAGAGCTG GTTGCGTGAG TCTCCTCAGC TCTGCTTACC GGTGCGACTA

51 _GCGGCAGCGA CGCGGCTAAA AGCGAAGGGG CGAGTGCGAG TCCCCTGAGC

LDL TGTACGAACG CGGTCGCCAT GGACCGCCCA GATGAGGGGC CTCCGGCCAA

L51 GACCCGCCGC CTGAGCAGCT CCGAGTCTCC ACAGCGCGAC CCGCCCCGC

2DL CGCCGCCGCC GCCGCCGCT CTCCGACTGC CGCTGCCTCC ACCCCAGCAG

251 CGCCGAGGC TCCAGGAGGA AACGGAGGCG GCACAGGTGC TGGCCGATAT

3DL GAGGGGGGTG GGACTGGGCC CCGCGCTGCC CCCGCCCT CCCTATGTCA

351 TTCTCGAGGA GGGGGGGATC CGCGCATACT TCACGCTCGG TGCTGAGTGT

4DL CCCGGCTGGG ATTCTACCAT CGAGTCGGG TATGGGGAAG CGCCCCCGCC

451 CACGGAGAGC CTGGAAGCAC TCCCCACTCC TGAGGCCTCG GGGGGGAGCC

5DL TGGAAATCGA TTTTCAGGTT GTACAGTCGA GCAGTTTTGG TGGAGAGGGG

551 GCCCTAGAAA CCTGTAGCGC AGTGGGGTGG GCGCCCCAGA GGTTAGTTGA

LDL CCCGAAGAGC AAGGAAGAGG CGATCATCAT AGTGGAGGAT GAGGATGAGG L AGGAGAGCTG GTTGCGTGAG TCTCCTCAGC TCTGCTTACC GGTGCGACTA 35 40 45 LOD CCCGAAGAGC AAGGAAGAGG CGATCATCAT AGTGGAGGAT GAGGATGAGG 701 AGGAAGCAGA GGAAGGTGAA GAGGGAAAGC AGAGAGAGAA ATGCCGAGAG 751 GATGGAGAC ATCCTGCAGG CACTGGAGGA TATTCAGCTG GATCTGGAGG
BOL CAGTGAACAT CAAGGCAGGC AAAGCCTTCC TGCGTCTCAA GCGCAAGTTC
B51 ATCCAGATGC GAAGACCCTT CCTGGAGCGC AGAGACCTCA TCATCCAGCA 50 PD1 TATCCCAGGC TTCTGGGTCA AAGCATTCCT CAACCACCCC AGAATTTCAA 951 TTTTGATCAA CCGACGTGAT GAAGACATTT TCCGCTACTT GACCAATCTG LODI CAGGTACAGG ATCTCAGACA TATCTCCATG GGCTACAAAA TGAAGCTGTA 55 1051 CTTCCAGACT AACCCCTACT TCACAAACAT GGTGATTGTC AAGGAGTTCC LIDI AGCGCAACCG CTCAGGCCGG CTGGTGTCTC ACTCAACCCC AATCCGCTGG 1151 CACCGGGGCC AGGAACCCCA GGCCCGTCGT CACGGGAACC AGGATGCGAG 1201 CCACAGCTTT TTCAGCTGGT TCTCAAACCA TAGCCTCCCA GAGGCTGACA

	€>	
.		

1251 GGATTGCTGA GATTATCAAG AATGATCTGT GGGTTAACCC TCTACGCTAC AAGTAAAGAA AAAGGGGCTC CAGATAAAG AGAAAGAC AAGAAAA AAGTAAAGT 1351 GAAACGTAAA ACCAGGGGCA GATGTGAGGT GGTGATCATG GAAGACGCCC 1401 CTGACTATTA TGCAGTGGAA GACATTTTCA GCGAGATCTC AGACATTGAT 1451 GAGACAATTC ATGACATCAA GATCTCTGAC TTCATGGAGA CCACCGACTA 5 1501 CTTCGAGACC ACTGACAATG AGATAACTGA CATCAATGAG AACATCTGCG 1551 ACAGCGAGAA TCCTGACCAC AATGAGGTCC CCAACAACGA GACCACTGAT 1601 AACAACGAGA GTGCTGATGA CCACGAAACC ACTGACAACA ATGAGAGTGC 1651 AGATGACAAC AACGAGAATC CTGAAGACAA TAACAAGAAC ACTGATGACA 1701 ACGAAGAGA CCCTAACAAC AACGAGAACA CTTACGGCAA CAACTTCTTC 10 1751 AAAGGTGGCT TCTGGGGCAG CCATGGCAAC AACCAGGACA GCAGCGACAG LBOL TGACAATGAA GCAGATGAGG CCAGTGATGA TGAAGATAAT GATGGCAACG LB5L AAGGTGACAA TGAGGGCAGT GATGATGATG GCAATGAAGG TGACAATGAA LADI CCCACCATC ALCACCACAC ACACALLEAC LACTALCACA AVELLALLE 1951 AGACTTTGAC AAGGATCAGG CTGACTACGA GGACGTGATA GAGATCATCT 15 20 25 2601 GAGGCGCTGC TGCCACCTTC CTCTCCCAAG TTCTTTCTCC ATCCCTCTCC 2651 TCTTCCCGCC GCGCCGCTAG CCCGCCTCGG TGTCTATGCA AGGCCGCTTC 2701 GCCATTGCGG TATTCTTTGC GGTATTCTTG TCCCCGTCCC CCAGAAGGCT 30 2751 CGCCTCTCCC CGTGGACCCT GTTAATCCCA ATAAAATTCT GAGCAAGTTT AAAAAAA AAAAAAA LOBS

35 BLAST Results

No BLAST result

40

Medline entries

98399864:
45 Vogel To Dittrich On Mehraein Yo Dechend Fo Schnieders Fo Schmidtke
Johnwine and human TSPYL genes: novel members of the
TSPY-SET-NAPILL family. Cytogenet Cell Genet 1998;81(3-4):265-70

50

Peptide information for frame 2

55

ORF from 119 bp to 2197 bp; peptide length: 693 Category: similarity to known protein Classification: unclassified

WO 01/98454 PCT/IB01/02050

```
1 MDRPDEGPPA KTRRLSSSES PQRDPPPPPP PPPLLRLPLP PPQQRPRLQE
       51 ETEAAQVLAD MRGVGLGPAL PPPPPYVILE EGGIRAYFTL GAECPGWDST
      DDD IESGYGEAPP PTESLEALPT PEASGSSLEI DFQVVQSSSF GGEGALETCS
      151 AVGWAPQRLV DPKSKEEAII IVEDEDEDER ESMRSSRRR RRRRRKQRKV
5
      201 KRESRERNAE RMESILQALE DIQLDLEAVN IKAGKAFLRL KRKFIQMRRP
      251 FLERRDLIIQ HIPGFWVKAF LNHPRISILI NRRDEDIFRY LTNLQVQDLR
      301 HISMGYKMKL YFQTNPYFTN MVIVKEFQRN RSGRLVSHST PIRWHRGQEP
      351 QARRHGNQDA SHSFFSWFSN HSLPEADRIA EIIKNDLWVN PLRYYLRERG
      401 SRIKRKKQEM KKRKTRGRCE VVIMEDAPDY YAVEDIFSEI SDIDETIHDI
10
      451 KISDFMETTD YFETTDNEIT DINENICDSE NPDHNEVPNN ETTDNNESAD
      501 DHETTDNNES ADDNNENPED NNKNTDDNEE NPNNNENTYG NNFFKGGFWG
      551 SHGNNQDSSD SDNEADEASD DEDNDGNEGD NEGSDDDGNE GDNEGSDDDD
      POP BDIEALER I EDEDKDGADA EDAIEIIZDE ZAEERIEER IGGDEDIAEE
15
      651 GNYEEEGSED VWEEGEDSDD SDLEDVLQVP NGWANPGKRG KTG
```

BLASTP hits

20

No BLASTP hits available

Alert BLASTP hits for DKFZphamy2_7j5, frame 2

- 25 TREMBL: ABO15345_1 gene: "HRIHFB2216"; Homo sapiens HRIHFB2216 mRNA, partial cds., N = 4, Score = 1393, P = 2.1e-165
- TREMBL:HSDJ486I3_2 gene: "dJ486I3.2"; product: "dJ486I3.2

 (KIAAD723

 (NAP (Nucleosome Assembly Protein) domain containg protein))";

 Human

DNA sequence from clone 486I3 on chromosome 6q22.1-22.3. Contains the

- 35 part of a gene for a novel protein, the gene for KIAAO721 (NAP (Nucleosome Assembly Protein) domain containg protein), the TSPYL gene for TSPY-like (testis specific protein, Y-linked like), and an RPSS
- 40 (40S Ribosomal Protein S5) pseudogene. Contains ESTs, STSs, GSSs and two putative CpG islands., N = 1, Score = 570, P = 3.4e-55
- 45 >TREMBL:ABO15345_1 gene: "HRIHFB2216"; Homo sapiens HRIHFB2216 mRNA.

partial cds. Length = 486

50 HSPs:

55

Score = 1393 (209.0 bits), Expect = 2.le-165, Sum P(4) = 2.le-165 165 Identities = 268/295 (90%), Positives = 268/295 (90%)

Query: 208
NAERMESILQALEDIQLDLEAVNIKAGKAFLRLKRKFIQMRRPFLERRDLIIQHIPGFWV 267

NAERMESILQALEDIQLDLEAVNIKAGKAFLRLKRKFIQMRRPFLERRDLIIQHIPGFWV Sbjct: 1

NAERMESILQALEDIQLDLEAVNIKAGKAFLRLKRKFIQMRRPFLERRDLIIQHIPGFWV 60

Query: 268 KAFLNHPRISILINRRDEDIFRYLTNLQVQDLRHISMGYKMKLYFQTNPYFTNMVIVKEF 327

KAFLNHPRISILINRRDEDIFRYLTNLQVQDLRHISMGYKMKLYFQTNPYFTNMVIVKEF

10 Sbjct: 61 KAFLNHPRISILINRRDEDIFRYLTNLQVQDLRHISMGYKMKLYFQTNPYFTNMVIVKEF 120

15 QRNRSGRLVSHSTPIR
LPEADRIAEIIKNDL
Sbjct: 121

5

QRNRSGRLVSHSTPIRWHRGQEPQARRHGNQDASHSFFSWFSNHSLPEADRIAEIIKNDL 180

25 WVNPLRYYLRERGSRIKRKKQEMKKRKTRGRCEVVIMEDAPDYYAVEDIFSEISDIDETI 240

Query: 448 HDIKISDFMETTDYFETTDNEITDINENICDSENPDHNEVPNNETTDNNESADDH 502

30 HDIKISDFMETTDYFETTDNEITDINENICDSENPDHNEVPNNETTDNNESADDH Sbjct: 241 HDIKISDFMETTDYFETTDNEITDINENICDSENPDHNEVPNNETTDNNESADDH 295

Score = 117 (17.6 bits), Expect = 9.0e-19, Sum P(4) = 9.0e-19

35 Identities = 32/77 (41%), Positives = 44/77 (57%)

Query: 42b
DAPDYYAVEDIFSEISDIDETIHDIKISDFMETTDYFETTDNEITDINENICDSENPDHN 485
+ DY+ D +EI+DI+E I D E D+ E +NE TD NE+

40 D E D+N
Sbjct: 250 ETTDYFETTD--NEITDINENICD----SENPDHNEVPNNETTDNNESADDHETTDNN 301

Query: 486 EVP--NNETT-DNNESADDH 502 45 E NNE DNN++ DD+ Sbjct: 302 ESADDNNANTDDN 321

> Score = 94 (14.1 bits), Expect = 2.le-165, Sum P(4) = 2.le-165 Identities = 16/16 (100%), Positives = 16/16 (100%)

50
Query: 678 QVPNGWANPGKRGKTG 693
QVPNGWANPGKRGKTG
Sbict: 471 QVPNGWANPGKRGKTG 486

55 Score = 90 (13.5 bits), Expect = 9.9e-16, Sum P(4) = 9.9e-16 Identities = 34/85 (40%), Positives = 45/85 (52%)

-206-

Query: 42b DAPDYYAVEDIFSEISDIDETIHDIKISDFME----TTDYFETTDN-EITDINENICDS 479

+ DY+ D +EI+DI+E I D + D E TTD E+ D+ E TD

NE+ D+

5 Sbjct: 250 ETTDYFETTD-NEITDINENICDSENPDHNEVPNNETTDNNESADDHETTDNNESADDN 307

Query: 480 -ENPDHN-----EVPNN-ETTDNN 496

ENP+ N E PNN E T N

10 Sbjct: 308 NENPEDNNKNTDDNEENPNNNENTYGN 334

Score = 87 (13.1 bits), Expect = 2.1e-165, Sum P(4) = 2.1e-165 Identities = 14/14 (100%), Positives = 14/14 (100%)

15 Query: 543 FFKGGFWGSHGNNQ 556

FFKGGFWGSHGNNQ

Sbjct: 336 FFKGGFWGSHGNNQ 349

Score = 85 (12.8 bits), Expect = 2.1e-165, Sum P(4) = 2.1e-165

20 Identities = 16/18 (88%), Positives = 17/18 (94%)

Query: 601 RDIEYYEKVIEDFDKDQA 618

RDIEYYEK IEDFD+DQA

Sbjct: 394 RDIEYYEKGIEDFDRDQA 411

25
Score = 60 (9.0 bits), Expect = 5.3e-03, Sum P(4) = 5.3e-03
Identities = 21/66 (31%), Positives = 33/66 (50%)

Query: 42b DAPDYYAVEDIFSEISDIDETIHD-IKIS-

30 DFMETTDYFETTDNEITDINENICDSENPD 483

D DY V +I Z+ Z +E I + I + D E +Y E ++ + E+

DZ+ D

Sbjct: 409

DQADYEDVIEIISDESVEEEGIEEGIQQDEDIYEEGNYEEEGSEDVWEEGEDSDDSDLED 468

35

Query: 484 HNEVPN 489

+VPN

Sbict: 469 VLQVPN 474

40 Score = 49 (7.4 bits), Expect = 1.4e-06, Sum P(4) = 1.4e-06 Identities = 12/35 (34%), Positives = 21/35 (60%)

Query: 463 ETTDNEITDINENICDSENPDHNEVPNNETTDNNE 497
E +D+E D NE + + D NE +NE +D+++

45 Sbjct: 360 EASDDEDNDGNEGDNEGSDDDGNE-GDNEGSDDD 393

Score = 42 (6.3 bits), Expect = 7.2e-06, Sum P(4) = 7.2e-06 Identities = 11/37 (29%), Positives = 18/37 (48%)

50 Query: 465 TDNEITDINENICDSENPDHNEVPNNETTDNNESADD 501 +DNE + D E+ D NE N + D+ D+

Sbjct: 354 SDNEADEAS----DDEDNDGNEGDNEGSDDDGNEGDN 386

55 Pedant information for DKFZphamy2_7j5, frame 2

Report for DKFZphamy2_7j5-2

```
ELENGTHD 693
           79435.07
   EMMI
           4.45
 5
   [[q]
               TREMBL: ABD15345_1 gene: "HRIHFB2216"; Homo sapiens
   EHOMOL
   HRIHFB2216 mRNA, partial cds. Le-171
   IFUNCATI Ob-10 assembly of protein complexes
                                        ES. cerevisiae,
   YKR048c3 4e-05
10
   [FUNCAT] 03.22 cell cycle control and mitosis
                                        ES. cerevisiae,
   YKR048c3 4e-05
   EFUNCATD 03-04 budding, cell polarity and filament formation
       ES. cerevisiae, YKRO48cl 4e-D5
   EFUNCATI 09.13 biogenesis of chromosome structure
   cerevisiae, YKRO48cJ 4e-05
15
    EFUNCATD 30.10 nuclear organization ES. cerevisiae, YKRO48cD
   4e-05
   EBLOCKZI BPD5P4PH
   EBFOCK23
          BP02646E
   EBLOCK23
          PF00424A
20
   [BLOCK2]
          BLOO415N Synapsins proteins
   [BLOCKZ]
          BP02799E
   [BF0CKZ]
           BL00048 Protamine Pl proteins
   EBFOCKZ
          PR00049D
25
   EBLOCK23
          PF00956D
   EBFOCK23
          PF00956C
   EBLOCKSI PFOO9568
   [PIRKW]
              nucleus 8e-33
              phosphoprotein &e-33
   [PIRKW]
30
              alternative splicing Be-33
   [PIRKW]
   Alpha_Beta
   EKWD
           LOW_COMPLEXITY
                       35.35 %
   SEQ -- MDRPDEGPPAKTRRLSSSESPQRDPPPPPPPPLLRLPLPPPPQQRPRLQEETEAAQVLAD
-- 35
   SEG
       PRD
       MRGVGLGPALPPPPYVILEEGGIRAYFTLGAECPGWDSTIESGYGEAPPPTESLEALPT
   SEQ
40
   SEG
       ····XXXXXXXXXX
   PRD
       PEASGSLEIDFQVVQSSSFGGEGALETCSAVGWAPQRLVDPKSKEEAIIIVEDEDEDER
   SEQ
   SEG
       -----xxxxxxx
45
   PRD
       SEQ
       ESMRSSRRRRRRRKQRKVKRESRERNAERMESILQALEDIQLDLEAVNIKAGKAFLRL
   SEG
       PRD
       րերերեր անական անակա
50
       KRKFIQMRRPFLERRDLIIQHIPGFWVKAFLNHPRISILINRRDEDIFRYLTNLQVQDLR
   SEQ
   SEG
   PRD
       55
   SEQ
       HISMGYKMKLYFQTNPYFTNMVIVKEFQRNRSGRLVSHSTPIRWHRGQEPQARRHGNQDA
   SEG
       PRD
```

	w	O 01/98454 PCT/IB01/02050				
	SEQ SEG PRD	SHSFFSWFSNHSLPEADRIAEIIKNDLWVNPLRYYLRERGSRIKRKK@EMKKRKTRGRCExxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx				
5	SEQ SEG PRD	<pre>vvimeDapDyyaveDifSeiSDidetiHDikiSDfmettDyfettDneitDinenicDSe</pre>				
10	SEQ SEG PRD	NPDHNEVPNNETTDNNESADDHETTDNNESADDNNENPEDNNKNTDDNEENPNNNENTYG xxxxxxxxxxxxxxxxxxxxxxxxxxxxxx				
15	SEQ SEG PRD	NNELKERENGZHENNGDZZDZDNE V DE VZD DE DND E O DO D				
20	SEQ SEG PRD	RDIEYYEKVIEDFDKD@ADYEDVIEIISDESVEEEGIEEGI@@DEDIYEEGNYEEEGSEDxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx				
20	SEQ SEG PRD	VWEEGEDSDDSDLEDVLQVPNGWANPGKRGKTG xxxxxxxxxxxxxxxx				
25	(No Prosite data available for DKFZphamy2_7j5-2)					
	(No	Pfam data available for DKFZphamy2_7j5.2)				
30		Pedant information for DKFZphamy2_7j5, frame 3				
35	Report for DKFZphamy2_7j5-3					
40	EMWI EpIl	75.88				
40	EKM]	All_Alpha				
45	SEQ SEG PRD	MRTSATARILTTMRSPTTRPLITTRVLMTTKPLTTMRVQMTTTRILKTITRTLMTTKRTLxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx				
50	SEQ SEG PRD	TTTRTLTATTSSKVASGAAMATTRTAATVTMKQMRPVMMKIMMATKVTMRAVMMMAMKVT xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx				
55	SEQ SEG PRD	MKAAMMTTETLSTMRKLLKTLTRIRLTTRT xxxxxxxx.xxxxxxxxxxxxxxxxxxxxxxxxxx				

(No Prosite data available for DKFZphamy2_7j5.3)

(No Pfam data available for DKFZphamy2_7j5.3)

DKFZphfbr2_78cl2

group: nucleic acid management

10 DKFZphfbr2_78cl2 encodes a novel 528 amino acid protein with high csimilarity to glutamyl-tRNA (6ln) amidotransferase subunit A of the hyperthermophilic bacterium Aquifex aeolicus.

The novel protein contains one ATP/GTP-binding site motif A (P15 loop). This loop interacts with one of the phosphate groups of a
A or G nucleotide. It is found in numerous ATP- or GTP-binding
proteins, such as ATP synthase alpha and beta subunits, Myosin
heavy chains, Kinesin heavy chains and kinesin-like proteins,
Dynamins and dynamin-like proteins, several kinases, DNA and RNA
20 helicases, GTP-binding elongation factors and the Ras family of
GTP-binding proteins. The protein seems to be expressed
ubiquitously.

The new protein can find application in the modulation of translational pathways.

similarity to glutamyl-tRNA (Gln) amidotransferase subunit A
 (Aquifex
30 aeolicus)

Sequenced by MediGenomix

35

Locus: /map="b8b.3 cR from top of Chrb linkage group"

Insert length: 3244 bp
Poly A stretch at pos. 3222, polyadenylation signal at pos. 3204

40

1 AGTGACAATT AAAGATGGCT GCGCCCATGT AACATCACTA GCGACCGGTG
51 ACCTCTTTTT CCCCCCTTGCC TGGCTCCTGT GGTGGCAGGC TGGGCACGAG
101 GACCATGCTG GGCCGGAGCC TCCGAGAAGT TTCTGCGGCA CTGAAACAAG
151 GCCAAATTAC ACCAACAGAG CTCTGTCAAA AATGTCTCTC TCTTATCAAG
201 AAGGCCAAGT TTCTAAAATGC CTACATTACT GTGTCAGAAG AGGTGGCCTT
251 AAAACAAGCT GAAGAATCAG AAAAGAGATA TAAGAATGGA CAGTCACTTG
301 GGGATTTAGA TGGAATTCCT ATTGCAGTAA AAGACAATTT CAGCACTTCT
351 GGCATTGAGA CAACATGTGC ATCAAATATG CTGAAAGGTT ATATACCACC
401 TTATAATGCT ACAGTAGTTC AGAAGTTGTT GGATCAGGGA GCTCTACTAA
451 TGGGAAAAAC AAATTTAGAT GAGTTTGCTA TGGGATCTGG GAGCACAGAT
50 501 GGTGTATTTG GACCAGTTAA AAACCCCTGG AGTTATTCAA AACGATATAG
551 AGAAAAGAGG AAGCAGAATC CCCCACAGCGA GAATGAAGAT TCAGACTGGC
601 TGATAACTGG AGGAAGCCCA GGTGGGAGTG CAGCTGCTGT ATCGGCGTTC
651 ACATGCTACG CGGCTTTAGG ATCAGATACA GGAGGATCGA CCAGAAATCC
701 TGCTGCCCAC TGTGGGCTTG TTGGTTTCAA ACCAAGCTAT GGCTTAGTTT
55 751 CCCGTCATGG TCTCATTCCC CTGGTGAAATT CGATGGTG CACTGGCCGG
B51 ACCTGACCCC AGGGACTCTA CCACAGTACA TGAACCTATT AATAAACCAT
901 TCATGCTTCC CAGTTTGGCA GATGTGAGCA AACTATGTAT AGGAATTCCA

						- 01.1201.02000
	751	AAGGAATATC	TTGTACCGGA	ATTATCAAGT	GAAGTACAGT	CTCTTTGGTC
	1001	CAAAGCTGCT	GACCTCTTTG	AGTCTGAGGG	GGCCAAAGTA	ATTGAAGTAT
	1051	CCCTTCCTCA	CACCAGTTAT	TCAATTGTCT	GCTACCATGT	ATTGTGCACA
	1101	TCAGAAGTGG	CATCGAATAT	GGCAAGATTT	GATGGGCTAC	AATATGGTCA
5	1151	CAGATGTGAC	ATTGATGTGT	CCACTGAAGC	CATGTATGCT	GCAACCAGAC
3	7507	GAGAAGGATT	TAATGATGTG	GTGAGAGGAA	GAATTCTCTC	AGGAAACTTT
	7527	TTCTTATTAA			TTTGTCAAAG	
			AAGAAAACTA	TGAAAATTAT		CACAGAAAGT
	7307	GAGACGCCTC	ATTGCTAATG	ACTTTGTAAA	TGCTTTTAAC	TCTGGAGTAG
	1351	ATGTCTTGCT	AACTCCCACC	ACCTTGAGTG	AGGCAGTACC	ATACTTGGAG
10	1401	TTCATCAAAG	AGGACAACAG	AACCCGAAGT	GCCCAGGATG	ATATTTTTAC
	1451	ACAAGCTGTA	AATATGGCAG	GATTGCCAGC	AGTGAGTATC	CCTGTTGCAC
	1501	TCTCAAACCA	AGGGTTGCCA	ATAGGACTGC	AGTTTATTGG	ACGTGCGTTT
	1551	TGTGACCAGC	AGCTTCTTAC	AGTAGCCAAA	TGGTTTGAAA	AACAAGTACA
	7P07	GTTTCCTGTT	ATTCAACTTC	AAGAACTCAT	GGATGATTGT	TCAGCAGTCC
15	1651	TTGAAAATGA	AAAGTTAGCC	TCTGTCTCTC	TAAAACAGTA	AACATATCTT
	1701	ACAAATTAAA	ATGACTTTTA	GGCTGGGTGC	AGTGGCTCAC	ACCTGTAATC
	1751	CCAGCACTTT	GGGAGGCCAA	GGCGAGCGGA	TCATGAGGTC	AGAAGATCTA
	1801	GAACAGCCTG	GTCAACATGG	TGAAACCCCG	TCTCTACTAA	AAATACAAAA
	1851	ATTAGCCAGG	CTTAGTGGCG	GGCATCTGTA	GTCCCAGCTA	CTCAGGAGGC
20	1901	TGAGGCAGGA	GAATCACTTG	AACCCTGGAG	GTGGAGGTTG	CAGTGAGCCG
20	1951	AGATCATGCC	ACTGCACTGC	ACTCCAGCCT	GGGTGACAAA	GCAAGACTGT
	5007	GTCTCAAAAT	AATAAATAA	AATAAAATAA	AATGACGTAC	AGAGATTCTA
	2051	TATTCTAGAG	AGTCAAATGG	TCTTGCTCAA	TTCTTGTAAT	TAGGTTCTTG
	5707	TTAATACAGT	CATTCCATGG	AATTACTTTT	TAAAATTCCT	GTGACAATTA
25	2727	ATAATACAGT	ACGTGTCAGC	ATTTAGTAAG	CATCCACTAA	GTGTACAATA
23	5507	CTTCTACAAT	AACACAAGAT	ACCTGTTCCT	CAAAGACAAT	
	2251	ATAATGTTCA	TTAAAGAGTT	TACAGTAAAA	ATAAGATTAG	GCATTCTGCC
	5307	CTCAAAAATT	GTACATCTGT			GGATAAACTT
	2351	GTCCTTCTAG	AGGTAACTTG	GTAACTAAAG	CACTAACAAA	AACATGAATA
20				GATAGCCTAG	GCAGGCAACT	TATCATGTGG
30	2401	TGAAGGCCGC	CTCAGGGGTT	GTTAAAAATG	CACAGAAACA	ATTGAGTGCG
	2451	ATTATTGGCT	TCTGAGCGCT	GAGCAGAGCA	GGTGGAAGAG	GAACTTTGAG
	2501	CACAGGAGGA	AATGCAACCA	GTCAGGGCCC	AGAATCATGC	AAATCTCAGG
	2551	GGTATGCCTC	TCTGGGGAGG	AGCTCCACTT	GCAGGGACTC	CTTTTATTTC
25	5P07	CCTAAGAAAG	AGCTGAAATG	ACTGAGAACT	TTCCTTTCCT	CCTTAGAGTT
35	5627	ACAATTTTAC	TTCTGCTATT	CCGGAGCCCA	TGCCTAGAAG	CCAGAACAAC
	2701	TCCATGTTAC	ACTGAGTTCA	TGCTCCTATT	TACTGATCAC	AAATGAGCTC
	2751	ATTAATGTCA	TCGAAACATT	TATTGTAACC	TAACAGACCA	TCACAGATTG
	5907	GAAACTTGGT	AGATAGCAGA	GCATGGTATT	AGTGAAAAAG	GTTCAAAATA
	2851	CACAAGTAAC	ATACACTCTG	AAAAACATGC	AGATAATTTG	CTGATGAAGC
40	2901	AGAAGAGGGG	ATGCGCATGG	CAAGAACTTG	CCTTACCCCA	GATTCTCTAT
	2951	ATCTCATGGT	TTCCTTTTCC	TCTTGACTGT	CTTTACGAGT	GTTTTTTATT
	3007	TGGGACCCTC	GAGCCCAGAG	ATATTAATGG	ATATCTGTAT	TCAATATTT <i>G</i>
	3051	ACAAAATCTA	ATGGAAACCA	TCCATTTACT	CATGATAAGG	CTTCATCACT
	3707	GGATTTCTGT	GTCTTCACTA	GAACACCATT	GTCATCTCAT	ATTGATCAGG
45	3151	TATTTTAATC	TAGCACTTAC	ATATTGTTGA	TAAATGAAAG	CTGAATTGTT
	3507	ACTTAATAAA	TTCACTTTGT	TTAGCAAAAA	AAAAAAAAA	AAAA

50 BLAST Results

No BLAST result

55 Medline entries

No Medline entry

Peptide information for frame 3

5 ORF from 105 bp to 1688 bp; peptide length: 528 Category: similarity to known protein Classification: Protein management Prosite motifs: ATP_GTP_A (112-119) 10 1 MLGRSLREVS AALK@G@ITP TELC@KCLSL IKKAKFLNAY ITVSEEVALK 51 QAEESEKRYK NGQSLGDLDG IPIAVKDNFS TSGIETTCAS NMLKGYIPPY 101 NATVVQKLLD QGALLMGKTN LDEFAMGSGS TDGVFGPVKN PWSYSKRYRE 15 151 KRKQNPHSEN EDSDWLITGG SPGGSAAAVS AFTCYAALGS DTGGSTRNPA 201 AHCGLVGFKP SYGLVSRHGL IPLVNSMDVP GILTRCVDDA AIVLGALAGP 251 DPRDSTTVHE PINKPFMLPS LADVSKLCIG IPKEYLVPEL SSEVQSLWSK 301 AADLFESEGA KVIEVSLPHT SYSIVCYHVL CTSEVASNMA RFDGLQYGHR 351 CDIDVSTEAM YAATREGFN DVVRGRILSG NFFLLKENYE NYFVKAQKVR 20 401 RLIANDFVNA FNSGVDVLLT PTTLSEAVPY LEFIKEDNRT RSAQDDIFTQ 451 AVNMAGLPAV SIPVALSNAG LPIGLAFIGR AFCDAALLTV AKWFEKAVAF 501 PVIQLQELMD DCSAVLENEK LASVSLKQ 25 BLASTP hits No BLASTP hits available 30 Alert BLASTP hits for DKFZphfbr2_78cl2, frame 3 PIR:F70322 glutamyl-tRNA (Gln) amidotransferase subunit A -Aquifex 35 aeolicus, N = 2, Score = 620, P = 4.3e-89>PIR:F70322 qlutamyl-tRNA (Gln) amidotransferase subunit A -Aquifex 40 aeolicus Length = 478 HSPs: 45 Score = 620 (93.0 bits), Expect = 4.3e-89, Sum P(2) = 4.3e-89 Identities = 135/319 (42%), Positives = 195/319 (61%) Querv: 187 ALGSDTGGSTRNPAAHCGLVGFKPSYGLVSRHGLIPLVNSMDVPGILTRCVDDAAIVLGA 246 50 +LGSDTGGS R PA+ CG++G KP+YG VSR+GL+ +S+D G+ R +D A+VL Sbict: 163 SLGSDTGGSIR@PASFCGVIGIKPTYGRVSRYGLVAFASSLD@IGVFGRRTEDVALVLEV 222 55 Query: LAGPDPRDSTTVHEPINKPFMLPSLADVSKLCIGIPKEYLVPELSSEVQSLWSKAADLFE 306 ++G D +DST+ P+ + +V L IG+PKE+ EL +V+ +

-34AASTSQMDEKDSTSAKVPVPE-

MZEEAKKEAKCTKICTLKELLEAETGLGAAKEVLENLIKETE 597

Query: 307

5 SEGAKVIEVSLPHTSYSIVCYHVLCTSEVASNMARFDGLQYGHRCDIDVSTEAMYAATRR 366 EG ++ EVSLPH YSI Y+++ SE +SN+AR+DG++YG+R

MYA TR

Sbjct: 282

KEGFEIKEVSLPHVKYSIPTYYIIAPSEASSNLARYDGVRYGYRAKEYKDIFEMYARTRD 341

10

20

Query: 367
EGFNDVVRGRILSGNFFLLKENYENYFVKAQKVRRLIANDFVNAFNSGVDVLLTPTTLSE 426
EGF V+ RI+ G F L Y+ Y++KAQKVRRLI NDF+ AF VDV-

+PTT

15 Sbjct: 342 EGFGPEVKRRIMLGTFALSAGYYDAYYLKAQKVRRLITNDFLKAFEE-VDVIASPT--P 398

Query: 427

AVPYLEFIKEDNRTRSAQDDIFTQAVNMAGLPAVSIPVALSNQGLPIGLQFIGRAFCDQQ 486 +P+ + M DI T N+AGLPA+SIP+A + GLP+G Q

TC+ + +

Sbjct: _399 TLPFKFGERLENPIEMYLSDILTVPANLAGLPAISIPIAWKD-GLPVGGQLIGKHWDETT 457

25 Query: 487 LLTVAK-WFEKQVQFPVIQL 505

LL ++ W +K + I L

Sbjct: 458 LLQISYLWEQKFKHYEKIPL 477

Score = 289 (43.4 bits), Expect = 4.3e-89, Sum P(2) = 4.3e-89 30 Identities = 64/143 (44%), Positives = 90/143 (62%)

Query: 4 RSLREVSAALKQGQITPTELCQKCLSLIKKAKF-LNAYITVSEEVALKQAEESEKRYKNG 62

+SL E+ LK+G+++P E+ + + + + AYIT ALKQAE

KSLSELRELLKRGEVSPKEVVESFYDRYNQTEEKVKAYITPLYGKALKQAESLKER---- 60

Querv: 63

40 QSLGĎLDGIPIAVKDNFSTSGIETTCASNMLKGYIPPYNATVVQKLLDQGALLMGKTNLD 122 L L GIPIAVKDN G +TTCAS +L+ ++ PY+ATV+++L

GAL++GKTNLD

Sbjct: bl -EL-

PLFGIPIAVKDNILVEGEKTTCASKILENFVAPYDATVIERLKKAGALIVGKTNLD 118

45

Query: 123 EFAMGSGSTDGVFGPVKNPWSYSK 146

EFAMGS + F P KNPW +

Sbjct: 119 EFAMGSSTEYSAFFPTKNPWDLER 142

50

Pedant information for DKFZphfbr2_78cl2, frame 3

Report for DKFZphfbr2_78cl2.3

55

ELENGTHI 528 EMWI 57468.78

```
[[q]
           5.57
   [HOMOL]
               PIR:E71725 glutamyl-tRNA amidotransferase chain A
   (gatA) RP152 - Rickettsia prowazekii 2e-93
   EFUNCATE r general function prediction
                                       [M. jannaschii,
   W1JJPOJ 96-PJ
   EFUNCATI 01.02.01 nitrogen and sulphur utilization
                                               EZ-
   cerevisiae, YMR293cJ Le-55
   EFUNCATE c energy conversion EM. gent
EFUNCATE Ol.Ol.10 amino-acid degradation
                              EM. genitalium, MGD991 4e-49
                                       ES. cerevisiae.
10
   YBR2O&cl 2e-3l
   EFUNCATI 01.03.01 purine-ribonucleotide metabolism
   cerevisiae, YBR20Acl 2e-31
   EBLOCKSI BLOO571
   ECI
           6.3.4.6 Urea carboxylase 5e-30
           3.5.1.4 Amidase 3e-39
15
   [EC]
           3.5.2.12 b-Aminohexanoate-cyclic-dimer hydrolase le-17
   ECI
   EPIRKUJ
               ligase 5e-3D
   EPIRKWI
               transmembrane protein 5e-30
   EPIRKWI
               ATP 5e-30
20
   EPIRKU
               crown gall tumor le-29
               mitochondrion 2e-13
   EPIRKUD
   EPIRKWI
               purine nucleotide binding 5e-30
               P-loop 5e-30
   EPIRKWI
   EPIRKU
               hydrolase 3e-39
25
               biotin 5e-3D
   EPIRKU3
           amidase 3e-39
   ESUPFAMI
   ESUPFAMI
           biotin carboxylase homology 5e-30
   ESUPFAMI
           indoleacetamide hydrolase 7e-92
   ESUPFAM3
           lipoyl/biotin-binding homology 5e-30
   EPROSITED ATP_GTP_A 1
30
   [KW]
           Alpha_Beta
           LOW_COMPLEXITY
   EKWI
                          2.46 %
35
   SEQ
       MLGRSLREVSAALK@G@ITPTELC@KCLSLIKKAKFLNAYITVSEEVALK@AEESEKRYK
   SEG
       PRD
       NGQSLGDLDGIPIAVKDNFSTSGIETTCASNMLKGYIPPYNATVVQKLLDQGALLMGKTN
   SEQ
40
   SEG
   PRD
       LDEFAMGZGZTDGVFGPVKNPWZYZKRYREKRKQNPHSENEDZDWLITGGZPGGZAAAVS
   SEQ
       -----xxxxxxxxxxx
   SEG
45
   PRD
       SEQ
       AFTCYAALGSDTGGSTRNPAAHCGLVGFKPSYGLVSRHGLIPLVNSMDVPGILTRCVDDA
   SEG
       PRD
       50
   ZEQ
       AIVLGALAGPDPRDSTTVHEPINKPFMLPSLADVSKLCIGIPKEYLVPELSSEV@SLWSK
   SEG
   PRD
       SEQ
55
       AADLFESEGAKVIEVSLPHTSYSIVCYHVLCTSEVASNMARFDGLQYGHRCDIDVSTEAM
   SEG
       PRD
```

WO 01/98454 PCT/IB01/02050 YAATRREGFNDVVRGRILSGNFFLLKENYENYFVKAQKVRRLIANDFVNAFNSGVDVLLT SEQ SEG PRD PTTLSEAVPYLEFIKEDNRTRSAQDDIFTQAVNMAGLPAVSIPVALSNQGLPIGLQFIGR SEQ SEG PRD SEQ AFCDQQLLTVAKWFEKQVQFPVIQLQELMDDCSAVLENEKLASVSLKQ 10 SEG PRD ccchhhhhhhhhhhhhhhhheeehhhhhhheeeecccceeeeccc 15 Prosite for DKFZphfbr2_78cl2.3 **PZ00073** 115->150 ATP_GTP_A PDOCUULT (No Pfam data available for DKFZphfbr2_7&cl2.3) 20 DKFZphfbr2_78dl8 25 group: brain derived DKFZphfbr2_78dl8 encodes a novel 535 amino acid protein with weak similarity to a human putative mitogen-activated protein kinase 30 kinase kinase. No informative BLAST results; No predictive prosite, pfam or SCOP motife-35 The new protein can find application in studying the expression profile of brain-specific genes. similarity to putative mitogen-activated protein kinase kinase kinase 40 (Homo sapiens) Sequenced by MediGenomix Locus: unknown 45 Insert length: 2158 bp Poly A stretch at pos. 2138, polyadenylation signal at pos. 2117 LATCCGGGGCC CCGGAACCCG AGCTGGAGCT GAAGCGCAGG CTGCGGGGCG

5L CGGAGTCGGG AGTGCAGGCC TGAGTGTTCC TTCCAGCATG TCGGAGGGGG

LOL AGTCCCAGAC AGTACTTAGC AGTGGCTCAG ACCCAAAGGT AGAATCCTCA

LSL TCTTCAGCCC CTGGCCTGAC ATCAGTGTCA CCTCCTGTGA CCTCCACAAC

COL CTCAGCTGCT TCCCCAGAGG AAGAAGAAGA AAGTGAAGAT GAGTCTGAGA

25L TTTTGGAAGA GTCGCCCTGT GGGCGCTGGC AGAAGAGGCG AGAAGAGGTG

BOL AATCAACGGA ATGTACCAGG TATTGACAGT GCATACCTGG CCATGGATAC

35L AGAGGAAGGT GTAGAGGTTG TGTGGAATGA GGTACAGTTC TCTGAACGCA

40L AGAACTACAA GCTGCAGGAG GAAAAGGTTC GTGCTGTTT TGATAATCTC 50 55

4D1 AGAACTACAA GCTGCAGGAG GAAAAGGTTC GTGCTGTGTT TGATAATCTG

WO 01/98454 PCT/IB01/02050 45% ATTCAATTGG AGCATCTTAA CATTGTTAAG TTTCACAAAT ATTGGGCTGA 5D1 CATTAAAGAG AACAAGGCCA GGGTCATTTT TATCACAGAA TACATGTCAT 551 CTGGGAGTCT GAAGCAATTT CTGAAGAAGA CCACAAGACG LOL ATGAATGAAA AGGCATGGAA GCGTTGGTGC ACACAAATCC TCTCTGCCCT
LSL AAGCTACCTG CACTCCTGTG ACCCCCCCAT CATCCATGGG AACCTGACCT
701 GTGACACCAT CTTCATCCAG CACAACGGAC TCATCAAGAT TGGCTCTGTG 5 751 GCTCCTGACA CTATCAACAA TCATGTGAAG ACTTGTCGAG AAGAGCAGAA
801 GAATCTACAC TTCTTTGCAC CAGAGTATGG AGAAGTCACT AATGTGACAA
851 CAGCAGTGGA CATCTACTCC TTTGGCATGT GTGCACTGGA GATGGCAGTG
901 CTGGAGATTC AGGGCAATGG AGAGTCCTCA TATGTGCCAC AGGAAGCCAT 10 951 CAGCAGTGCC ATCCAGCTTC TAGAAGACCC ATTACAGAGG GAGTTCATTC 1001 AAAAGTGCCT GCAGTCTGAG CCTGCTCGCA GACCAACAGC CAGAGAACTC 1051 CTGTTCCACC CAGCATTGTT TGAAGTGCCC TCGCTCAAAC TCCTTGCGGC LIDI CCACTGCATT GTGGGACACC AACACATGAT CCCAGAGAAC GCTCTAGAGG 1151 AGATCACCAA AAACATGGAT ACTAGTGCCG TACTGGCTGA AATCCCTGCA 15 1201 GGACCAGGAA GAGAACCAGT TCAGACTTTG TACTCTCAGT CACCAGCTCT 1251 GGAATTAGAT AAATTCCTTG AAGATGTCAG GAATGGGATC TATCCTCTGA LIBLE CAGCCTTTGG GCTGCCTCGG CCCCAGCAGC CACAGCAGGA GGAGGTGACA 1351 TCACCTGTCG TGCCCCCCTC TGTCAAGACT CCGACACCTG ACCCAGCTGA 1401 GGTGGAGACT CGCAAGGTGG TGCTGATGCA GTGCAACATT GAGTCGGTGG 20 1451 AGGAGGAGT CAAACACCAC CTGACACTTC TGCTGAAGTT GGAGGACAAA 1501 CTGAACCGGC ACCTGAGCTG TGACCTGATG CCAAATGAGA ATATCCCCGA 1551 GTTGGCGGCT GAGCTGGTGC AGCTGGGCTT CATTAGTGAG GCTGACCAGA ILOI GCCGGTTGAC TTCTCTGCTA GAAGAGACCT TGAACAAGTT CAATTTTGCC 1651 AGGAACAGTA CCCTCAACTC AGCCGCTGTC ACCGTCTCCT CTTAGAGCTC 25 1701 ACTCGGGCCA GGCCCTGATC TGCGCTGTGG CTGTCCCTGG ACGTGCTGCA 1751 GCCCTCCTGT CCCTTCCCC CAGTCAGTAT TACCCTGTGA AGCCCCTTCC
1801 CTCCTTTATT ATTCAGGAGG GCTGGGGGGG CTCCCTGGTT CTGAGCATCA
1851 TCCTTTCCC TCCCCTCTT TCCTCCCCTC TGCACTTTGT TTACTTGTTT
1901 TGCACAGACG TGGGCCTGGG CCTTCTCAGC AGCCGCCTTC TAGTTTGGGGG
1951 CTAGTCGCTG ATCTGCCGGC TCCCGCCCAG CCTGTGTGGA AAGGAGGCCC 30 35 2351 AAAAAAA

BLAST Results

40 No BLAST result

45

Medline entries

No Mediine entry

50
Peptide information for frame 1

ORF from && bp to 1692 bp; peptide length: 535
55 Category: similarity to unknown protein
Classification: Protein management

1 MSEGESQTVL SSGSDPKVES SSSAPGLTSV SPPVTSTTSA ASPEEEEESE

51 DESEILEESP CGRWQKRREE VNQRNVPGID SAYLAMDTEE GVEVVWNEVQ
101 FSERKNYKLQ EEKVRAVFDN LIQLEHLNIV KFHKYWADIK ENKARVIFIT
151 EYMSSGSLKQ FLKKTKKNHK TMNEKAWKRW CTQILSALSY LHSCDPPIH
201 GNLTCDTIFI QHNGLIKIGS VAPDTINNHV KTCREEQKNL HFFAPEYGEV
5 251 TNVTTAVDIY SFGMCALEMA VLEIQGNGES SYVPQEAISS AIQLLEDPLQ
301 REFIQKCLQS EPARRPTARE LLFHPALFEV PSLKLLAAHC IVGHQHMIPE
351 NALEEITKNM DTSAVLAEIP AGPGREPVQT LYSQSPALEL DKFLEDVRNG
401 IYPLTAFGLP RPQQPQQEEV TSPVVPPSVK TPTPEPAEVE TRKVVLMQCN
451 IESVEEGVKH HLTLLLKLED KLNRHLSCDL MPNENIPELA AELVQLGFIS
10 501 EADQSRLTSL LEETLNKFNF ARNSTLNSAA VTVSS

BLASTP hits

15

No BLASTP hits available

Alert BLASTP hits for DKFZphfbr2_78dl8, frame 1

20 TREMBL:ACOO9465_14 gene: "T9J14.14"; product: "putative mitogen activated protein kinase kinase"; Arabidopsis thaliana chromosome III

BAC T9J14 genomic sequence; complete sequence; N = 1, Score = 372, P =

25 1-9e-33

TREMBL:AF145690_1 gene: "BcDNA.LD2865?"; product: "BcDNA.LD2865?"; Drosophila melanogaster clone LD28657 BcDNA.LD2865?

30 (BcDNA-LD28657)

mRNA, complete cds., N = 1, Score = 1140, P = 1.3e-115

PIR:TO2951 probable mitogen activated protein kinase - rice, N = 1,

35 Score = 391_{1} P = 1.4e-35

>TREMBL:AF145690_1 gene: "BcDNA.LD2865?"; product: "BcDNA.LD2865?";

40 Drosophila melanogaster clone LD28657 BcDNA-LD28657 (BcDNA-LD28657) mRNAcomplete cds-

Length = 637

45 HSPs:

Score = 1140 (171.0 bits) Expect = 1.3e-115 P = 1.3e-115
Identities = 230/465 (49%) Positives = 304/465 (65%)

50 Query: 61
CGRWQKRREEVNQRNVPGIDSAYLAMDTEEGVEVVWNEVQFSERKNYKLQEEKVRAVFDN 120
CGRW KRREEV+QR+VPGID +LAMDTEEGVEVVWNEVQ++ + K
QEEK+R VFDN

Sbict: 102

55 CGRWLKRREEVDQRDVPGIDCVHLAMDTEEGVEVVWNEVQYASLQELKSQEEKMRQVFDN 161

Query: 121 LIQLEHLNIVKFHKYWADIKE-NKARVIFITEYMSSGSLKQFLKKTKKNHKTMNEKAWKR 179

L+QL+H NIVKFH+YW D ++ + RV+FITEYMSSGSLKQFLK+TK+N K

+ ++W+R

Sbjct: 162

LLQLDHQNIVKFHRYWTDTQQAERPRVVFITEYMSSGSLKQFLKRTKRNAKRLPLESWRR 221

5

Query: 180

WCTQILSALSYLHCOTIFICATIONHOLIKIGSVAPDIIHGHCOZHLYCALSZIDTDW ++ V ++ VQ VZDIXHCOTIFIQHCOZHLYCALSZIDTDW ++ V ++ VQ VZDIXHCOTIFIQHCOZHLYCALSZIDTDW

RE ++

10 Sbjct: 222

WCTQILSALSYLHSCSPPIIHGNLTCDSIFIQHNGLVKIGSVVPDAVHYSVRRGRERERE 281

Query: 24D ----LHFF-APEYGEVTVVTTAVDIYSFGMCALEMAVLEIQ-GNGESSYVPQEAISSAIQ 293

15

H+F APEYG +T A+DIY+FGMCALEMA LEIQ N ES+ +

+E I I

Sbjct: 282

RERGAHYFQAPEYGAADQLTAALDIYAFGMCALEMAALEIQPSNSESTAINEETIQRTIF 341

20 Query: 294 LLEDPLQREFIQKCLQSEPARRPTARELLFHPALFEVPSLKLLAAHCIV--- GHQHMIPE 350

LE+ LQR+ I+KCL +P RP+A +LLFHP LFEV SLKLL AHC+V

++ M E

Sbjct: 342

25 SLENDLQRDLIRKCLNPQPQDRPSANDLLFHPLLFEVHSLKLLTAHCLVFSPANRTMFSE 401

Query: 351 NALEEITKNMDTSAVLAEIPAGPGREPVQTLYSQSPALELDKFLEDVRNGIYPLTAFGL 409

30 G+YPL +

Sbjct: 402

TAFDGLMQRYYQPDVVMAQLRLAGGQERQYRLADVSGADKLEKFVEDVKYGVYPLITYS- 460

Querv: 410

Sbict: 461

GKKPPNFRSRAASPERADSVKSATPEPVDTESRRIVNMMCSVKIKEDSNDITMTILLRMD 520

40 Query: 470 XXXXXXXSCDLMPNENIPELAAELVQLGFISEADQSRLTSLLEETL 515

+C + N+ +L +ELV+LGF+ DQ ++ LLEETL

Sbjct: 521 DKMNRQLTCQVNENDTAADLTSELVRLGFVHLDDQDKIQVLLEETL 566

45

Pedant information for DKFZphfbr2_78dl8, frame 1

Report for DKFZphfbr2_78dl8-1

50

ELENGTHI 564

EMM3 62464.87

[pI] 5.10

EHOMOLI TREMBL: AF145690_1 gene: "BcDNA-LD28657"; product:

755 "BcDNA-LD28657"; Drosophila melanogaster clone LD28657
BcDNA-LD28657 (BcDNA-LD28657) mRNA, complete cds. le-123
EFUNCATI D3-22 cell cycle control and mitosis ES. cerevisiae,
YJLD95wl be-15

30.03 organization of cytoplasm EFUNCATI ES. cerevisiae -YJLD95wJ be-15

- [FUNCAT] ll.Ol stress response IS. cerevisiae, YJLO95wl be-l5 EFUNCATD 03.01 cell growth ES. cerevisiae, YJL095wl be-l5
- **EFUNCATI** 10.02.11 key kinases ES. cerevisiae, YJL095wl be-15 EFUNCATI D3.04 budding, cell polarity and filament formation
 ES. cerevisiae, YJLD95wl be-15
 - **EFUNCATI** 98 classification not yet clear-cut ES. cerevisiae; YLRO96wJ 2e-D9
- 10 **EFUNCATI** 30.02 organization of plasma membrane ES. cerevisiae; YLR096w3 2e-09 **EFUNCATI** 10.03.11 key kinases ES. cerevisiae, YNRD3lc3 3e-09

EFUNCATI 09.01 biogenesis of cell wall ES. cerevisiae.

YNR031c1 3e-09

- 03.07 pheromone response, mating-type determination, **EFUNCATI** IFUNCATO 10.05.11 key kinases IS. cerevisiae, YLR362w0 4e-08 **EFUNCATD** 10.04.11 key kinases ES- cerevisiae₁ YLR362w1 4e-08 11.04 dna repair (direct repair, base excision repair **EFUNCATI**
- 20 and nucleotide excision repair) ES. cerevisiae, YPL153cl le-07 EFUNCATI 03.19 recombination and dna repair ES. cerevisiae; YPL153cl le-07

EFUNCATI 03.22.01 cell cycle check point proteins EZcerevisiae, YPL153cl le-D7

- le-07 **EFUNCATI** 03.25 cytokinesis [S. cerevisiae, YDR507c] Le-07 **EFUNCATI** 10.99 other signal-transduction activities cerevisiae, YPL153cl le-07
- **EFUNCATI** O3.13 meiosis ES. cerevisiae, YDR523cl 3e-07 **EFUNCATE** 03.10 sporulation and germination ES. cerevisiae. YDR523c3 3e-07 **EFUNCATI** 03.16 dna synthesis and replication ES. cerevisiae; YMRDDlcl 2e-Ob
- EFUNCATI 99 unclassified proteins ES. cerevisiae, YDR490cl 3e-05 **EFUNCATE 05.07** translational control ES. cerevisiae, YDR283cl le-04

EFUNCATE 01.05.04 regulation of carbohydrate utilization

40 cerevisiae, YDR477wl le-04

> EBLOCKS] PFOOL37A BP03191J **EBFOCK21**

> **CBLOCKS3** PFO1317B

EZCOPI dlir3a_ 5.1.1.2.6 insulin receptor Complex

(transferase/substrate) 2e-53

dlphk___ 5.1.1.1.6 gamma-subunit of glycogen [CODI phosphorylase kinas 3e-68 ESCOP1 dlfgkb_ 5.1.1.2.5 Fibroblast growth factor receptor 1 Thuman (Hom le-55

50 [SCOP] dlabo___ 5.1.1.1.14 Protein kiase CK2, alpha subunit EMaize (Ze 2e-55 ESCOPI d3lck__ 5.1.1.2.2 Lymphocyte kinase (lck) (Homo sapiens) 7e-54 EZC0P] d2erk___ 5.1.1.1.1 MAP kinase Erk2 Erat (Rattus

9e-71 55 norvegicus)

dlcdkb_ 5.1.1.2 cAMP-dependent PK, catalytic ESCOPI subunit Comple Le-55

WO 01/98454 PCT/IB01/02050 dlhcl_ _ 5.l.l.l.l Cyclin-dependent PK [Human (Homo sapiens) 4e-67 2.7.1.112 Protein-tyrosine kinase 4e-0b [EC] 2.7.1.37 Protein kinase 3e-09 [[EC]] **EPIRKWI** phosphotransferase 2e-28 **EPIRKWI** nucleus 3e-06 [PIRKW] RNA binding 3e-10 tandem repeat 4e-07 **EPIRKU**J cell cycle control 3e-Ob [PIRKW] serine/threonine-specific protein kinase 2e-13 **EPIRKUJ** 10 transmembrane protein 4e-07 **EPIRKU**I [PIRKW] autophosphorylation 3e-10 [PIRKW] tyrosine-specific protein kinase 4e-0b [PIRKW] magnesium 4e-07 EL-es PTA 15 [PIRKW] **EPIRKWI** receptor 4e-07 **EPIRKWI** phosphoprotein 2e-13 [PIRKW] apoptosis 3e-06 [PIRKW] glycoprotein 4e-07 **EPIRKWI** 20 protein kinase 2e-28 **CPIRKWI** signal transduction 2e-D8 [PIRKW] cell division le-ll calmodulin binding 3e-Db **EPIRKUD** ESUPFAMD protein kinase byr2 le-Ob
ESUPFAMD unassigned Ser/Thr or Tyr-specific protein kinases 2e-25 73 ESUPFAMD leucine-rich alpha-2-glycoprotein repeat homology 4e-07 ESUPFAMI double-stranded RNA-binding repeat homology 3e-10
ESUPFAMI SAM homology 1e-0b
ESUPFAMI death-associated protein kinase 3e-0b
ESUPFAMI repeat homology 3e-0b
ESUPFAMI protein kinase homology 2e-28
ESUPFAMI kinase-related transforming protein 2e-0b
ESUPFAMI protein kinase SPK1 3e-0b
ESUPFAMI protein kinase Xa21 4e-07
ESUPFAMI protein kinase TIK 3e-10 30 35 ESUPFAMI kinase interaction domain homology 3e-06 EPFAM3 Eukaryotic protein kinase domain All_Alpha 40 EKWI EKWI 3 D [KW] LOW_COMPLEXITY 36.49 % SEQ IRGPGTRAGAEAQAAGRGVGSAGLSVPSSMSEGESQTVLSSGSDPKVESSSSAPGLTSVS SEG 1kobA

50 SEQ PPVTSTTSAASPEEEEESEDESEILEESPCGRWQKRREEVNQRNVPGIDSAYLAMDTEEG lkobA

SEQ VEVVUNEVQFSERKNYKLQEEKVRAVFDNLIQLEHLNIVKFHKYWADIKENKARVIFITE 55 ·····CHHHHHHHHHHHHHHHTTTBTTBCCEE----EEEETTTEEEEEEC

5	SEQ YMSSGSLKQFLKKTKKNHKTMNEKAWKRWCTQILSALSYLHSCDPPIIHGNLTCDTIFIQ SEG
	SEQ HNGLIKIGSVAPDTINNHVKTCREEQKNLHFFAPEYGEVTNVTTAVDIYSFGMCALEMAV SEG
10	Ъкора ттссееессттттеестттееееетттааассиннынисссвсыныныныныныныныны
	SEQ LEIQGNGESSYVPQEAISSAIQLLEDPLQREFIQKCLQSEPARRPTARELLFHPALFEVP SEG
15	ССТТТТСССНИННИННННССССТТТНИННИНННТТТТТБББССССИНИНННТТТТ
20	SEQ SLKLLAAHCIVGHQHMIPENALEEITKNMDTSAVLAEIPAGPGREPVQTLYSQSPALELD SEG
,	SEQ KFLEDVRNGIYPLTAFGLPRPQQPQQEEVTSPVVPPSVKTPTPEPAEVETRKVVLMQCNI SEGxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
25	
	SEQ ESVEEGVKHHLTLLLKLEDKLNRHLSCDLMPNENIPELAAELVQLGFISEADQSRLTSLL SEGxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
30	
- 01.	SEQ EETLNKFNFARNSTLNSAAVTVSS SEG
35	
	(No Prosite data available for DKFZphfbr2_?&dl&.l)
40	Pfam for DKFZphfbr2_78dl8.l
	HMM_NAME Eukaryotic protein kinase domain
45	HMM *rLnHPNIIRFYDwFed···ddDHIYMIMEYMeGGDLFDYIrrng····p +L H NI++F ++ D + ++ +I+EYM G+L +++++ +
50	Query 152 QLEHLNIVKFHKYWADIKENKARVIFITEYMSSGSLKQFLKKTKKNHKT 200
	HMM MsEweIrfIMy@ILrGMeYLHSMgIIHRDLKPENILIDeNgqIKIcDF M+E+ +++ +@IL++++YLHS IIH L + I+I +NG
55	IKI+ MNEKAWKRWCTQILSALSYLHSCDPPIIHGNLTCDTIFIQHNGLIKIGSV 250

HMM

GLARqMnnYerMttfCGTPWYMMAPEVIImgnyYttkVDMWSFGCILWEM

++ N+ + + + APE + ++ TT+VD++SFG+

EM

5 Query 251 APDTINNHVKTCREEQKNLHFF-APEY-GEVTNVTTAVDIYSFGMCALEM 298

HMM

MTGepPFyddnMemImrIiqrfrrpfWpnCSeElyDFMrwCWnyDPekRP

10

P++RP

Query 299 A--VLEIQ-

GNGESSYVPQEAISSAIQLLEDPLQREFIQKCLQSEPARRP 345

15 HMM

TFr@ILnHPWF*

T+R++L HP +

Query

346 TARELLFHPAL

356

DKFZphfbr2_78d4

5

20

30

group: transmembrane protein

DKFZphfbr2_78d4 encodes a novel 188 amino acid protein without similarity to known proteins.

The novel protein contains 1 transmembrane region and a

10 Cytochrome c family heme-binding site.

No informative BLAST results: No predictive prosite; pfam or SCOP motife.

The new protein can find application in studying the expression profile of brain-specific genes and as a new marker for amygdala cells.

weak similarity to hypothetical protein of Arabidopsis thaliana

perhaps complete cds.
Pedant: TRANSMEMBRANE

Sequenced by MediGenomix

25 Locus: unknown

Insert length: 1547 bp
Poly A stretch at pos. 1527, polyadenylation signal at pos. 1508

1 TTGCCGCCGC CGCCACCCC GCCCAGGATG GCGGAAGTGG AGGCGCCGAC
51 GGCGGCCGAG ACGGACATGA AGCAATATCA AGGCTCCGGC GGCGTCGCCA



1351 TGGGTGGCCT GCCAGCACAG CCAGTGCCAT CAGGGAGCTG AAGGGGCTGT 1401 CCCCCACCTA ACTCCAGCTC CCCCTTCACG TTGTCACCAA GGCCCTGTGC 1451 CGCCCGCCTC GCCCCCTGC TCTGTGGATT CCTTTGGGAA GGGCTCCCTG 1501 GGCAGGACAA TAAAGAGTTT TGACTCCAAA AAAAAAAAA AAAAAAA

5

BLAST Results

10 Entry TO2616 from database PIR: hypothetical protein Tl9Ll8.12 - Arabidopsis thaliana Score = 229, P = 1.3e-17, identities = 57/161, positives = 78/1617 frame +1

15

Medline entries

20

No Medline entry

25

Peptide information for frame 1

ORF from 28 bp to 591 bp; peptide length: 188 Category: similarity to unknown protein 30 Classification: no clue Prosite motifs: CYTOCHROME_C (121-119)

1 MAEVEAPTAA ETDMKQYQGS GGVAMDVERS RFPYCVVWTP IPVLTWFFPI 51 IGHMGICTST GVIRDFAGPY FVSEDNMAFG KPAKYWKLDP AQVYASGPNA 101 WDTAVHDASE EYKHRMHNLC CDNCHSHVAL ALNLMRYNNS TNWNMVTLCF 35 151 FCLLYGKYVS VGAFVKTWLP FILLLGIILT VSLVFNLR

40

BLASTP hits

No BLASTP hits available

45 Alert BLASTP hits for DKFZphfbr2_78d4, frame 1

PIR:TO2616 hypothetical protein T19L18.12 - Arabidopsis thaliana, 2_{1} Score = $22b_{1}$ P = 4.5e-21

50

>PIR:T02616 hypothetical protein T19L18.12 - Arabidopsis thaliana Length = 267

55 HSPs:

> Score = 226 (33.9 bits), Expect = 4.5e-21, Sum P(2) = 4.5e-21 Identities = 52/132 (39%), Positives = 71/132 (53%)

Query: MDVERSRFPYCVVWTPIPVLTWFFPIIGHMGICTSTGVIRDFAGPYFVSEDNMAFGKPAK 84 +D ++Z+FP C+VUTP+PV++W P IGH+G+C GVI DFAG F++ D+ Sbjct: IDTKKSKFPCCIVWTPLPVVSWLAPFIGHIGLCREDGVILDFAGSNFINVDDFAFGPPAR 120 85 YWKLDPAQVYASGPNAWDTAVHDASEEYKHRMHNLC--CDNCHSHVALALNLMRYNNST- 141 10 Y + LD ++KH Z NG Н YN T Sbict: 121 YLQLDRTKCCLP-PNMGG---HTCKYGFKHTDFGTARTWDNALSSSTRSFEHKTYNIFTC 176 15 142 NWN-MVTLCFFCLLYG 156 Query: N + V C LYG 177 NCHSFVANCLNRLCYG 192 Sbjct: 20 Score = 157 (23.6 bits), Expect = 1.8e-13, Sum P(2) = 1.8e-13 Identities = 27/81 (33%), Positives = 50/81 (61%) 707 WDTAVHDASEEYKHRMHNLCCDNCHSHVALALNLMRYNNSTNWNMVTLCFFCLLYGKYVS 160 25 WD A+ ++ ++H+ +N+ NCHS VA LN + Y S WNMV + ++ GK+++ Sbjct: 155 WDNALSSSTRSFEHKTYNIFTCNCHSFVANCLNRLCYGGSMEWNMVNVAILLMIKGKWIN 214 30 161 VGAFVKTWLPFILL--LGIIL 179 Query: + V+++LP: ++ LG++L 215 GSSVVRSFLPCAVVTSLGVVL 235 Sbict: Score = 36 (5.4 bits), Expect = 4.5e-21, Sum P(2) = 4.5e-21 35 Identities = 7/21 (33%), Positives = 14/21 (66%) 10 AETDMKQYQGSGGVAMDVERS 30 Query: ++ ++K +G G MD++RS Sbjct: 15 SDRNLKMZRGRGVPMMDLKRZ 35 40 Pedant information for DKFZphfbr2_78d4, frame 1 45 Report for DKFZphfbr2_78d4-1 CLENGTHD 188 EWWI 21178-66 50 6.27 [[q] PIR:T02616 hypothetical protein T19L18.12 -CHOMOLI Arabidopsis thaliana 7e-32 EPROSITE CYTOCHROME_C [KW] TRANSMEMBRANE 55

SEQ MAEVEAPTAAETDMKQYQGSGGVAMDVERSRFPYCVVWTPIPVLTWFFPIIGHMGICTST PRD cccccchhhhhhhhhhcccccccccccccccccccceeecce

	W	O 01/98454	PCT/IB01/02050			
	MEM		• • • • • • • • • • • • • • • • • • • •			
5	SEQ PRD MEM	GVIRDFAGPYFVSEDNMAFGKPAKYWKLDPAQVYASGPNAWD eeeeccccccccccccccccccccccccccccccccc	cccccchhhhhhhhhee			
10	SEQ PRD MEM	CDNCHSHVALALNLMRYNNSTNWNMVTLCFFCLLYGKYVSVG ecccchhhhhhhhhhhccccccchhhhhhhhhhccceeeee	eeeeeccceeeceec			
	SEQ PRD MEM	VSLVFNLR ceeeeccc MMMMM				
15			,			
		Prosite for DKFZphfbr2_78d	4.1			
20	0029	1390 121->127 CYTOCHROME_C	PD0C00369			
	(No	Pfam data available for DKFZphfbr2_78d4.1)	<u> </u>			

DKFZphfbr2_78el8

5 group: brain derived

DKFZphfbr2_78el8 encodes a novel 307 amino acid protein without similarity to known proteins.

The mRNA is differentially polyadenylated. No informative BLAST results: No predictive prosite: pfam or SCOP motife.

The new protein can find application in studying the expression profile of brain-specific genes.

similarity to hypothetical protein of Arabidopsis thaliana

- 20 differential polyadenylation
 > 7 exons
 complete on human genomic clone 451821apperhaps complete cds.
- 25 Sequenced by MediGenomix

Locus: /map="144.50 cR from top of Chrb linkage group"

Insert length: 3096 bp
30 Poly A stretch at pos. 3075, polyadenylation signal at pos. 3047

L TGGTGAGTTC GGAGTAGAGA TGGCCGCGCT TGCACCGCTG CCCCCGCTCC
51 CCGCACAGCT CAAGAGCATA CAGCATCATC TGAGGACGGC TCAGGAGCAT
101 GACAAGCGAG ACCCTGTGGT GGCTTATTAC TGTCGTTTAT ACGCAATGCA
151 GACTGGAATG AAGATCGATA GTAAAACTCC TGAATGTCGC AAATTTTTAT
201 CAAAGTTAAT GGATCAGTTA GAAGCTCTAA AGAAGCAGTT GGGTGATAAT
251 GAAGCTATTA CTCAAGAAAT AGTGGGCTGT GCCCATTTGG AGAATTATGC
301 TTTGAAAATG TTTTTGTATG CAGACAATGA AGATCGTGCT GGACGATTTC
351 ACAAAAACAT GATCAAGTCC TTCTATACTG CAAGTCTTTT GATAGATGTC
401 ATAACAGTAT TTGGAGAACT CACTGATGAA AATGTGAAAC ACAGGAAGTA
451 TGCCAGATGG AAGGCCAACAT ACATCCATAA TTGTTTAAAG AATGGGGAGA
501 CTCCTCAAGC AGGCCCTGTT GGAATTGAAG AAGATAATGA TATTGAAGAA
551 AATGAAGATG CTGGAGCAGC CTCTCTGCCC ACTCAGCCAA CTCAGCCATC
601 ATCATCTTCA ACTTATGACC CAAGCAACAT GCCATCAGCC AACTATACTG
651 GAATACAGAT TCCTCCGGGT GCACACGCTC CAGCTAATAC ACCAGCAGAA
701 GTGCCTCACA GCACAGGTGT AGCAAGTAAT ACTATCCAAC CTACTCCACA
751 GACTATACCT GCCATTGATC CCGCACTTTT CAATACAATT TCCCAGGGGG
801 ATGTTCGTCT AACCCCAGAA GACTTTGCTA GAGCTCAGAA GTACTGCAAA 35 40 45 BOL ATGTTCGTCT AACCCCAGAA GACTTTGCTA GAGCTCAGAA GTACTGCAAA 851 TATGCTGGCA GTGCTTTGCA GTATGAAGAT GTAAGCACTG CTGTCCAGAA 50 POL TCTACAAAAG GCTCTCAAGT TACTGACGAC AGGCAGAGAA TGAAGCCTTT 95% GTATGACAGA CCCATGTATT TTTGGCATGA GGAACTAACA GTCCATTACT BOOD CTATCTTCAG CCTATCAGGA TCACAGTTTT AAGGAAGACT TGGTTTTGTT 1051 GAATATGACA ATGAAATCTG TGTGTATCAG ATTTTTATTG AAGCATTCAT 1101 CAGCAGCCTC AACCAGTTTT CATTGTCCAT TTACTAGATT CAATCGTCTC 55 1151 TGAGTATATA GGGCTGATGT TAGCAAGACC CTAAAAATGT CCATTGAACC **ጌ**ደዐጌ **CTGCTTCAAA AAATGAAAAC ACACCTCTAT AAAATGTGTA CTG**GGAATAA 1251 GCTTTGTATT TACATACATT AGGGGAATTT TTTAAAATCT GTAATGTTTG

WO 01/98454

PCT/IB01/02050

L301 GACAAACAGA TGATATTACT TTGCTATAAA ATTATAAATG TAACTTTTAA
L351 TAAAGATAGC CAGAATATTC TAAATTAGAA ATTACGTTTT TGTTTCCCTC 1401 AAGACATAAA ACAAATATAA ACATTCTAAA CTGCTGGATG AATCTGAAAA 1451 GACATTAAGT TCAAATTTTA ATTTATTCTC ATATTAAATA TAACTCCATT 1501 AAAAGTTTAA AATTTCATGG GAGAAAATAT AATAAGGTAA AGAGGTAGAA 5 1551 TCACTTTCAG ACTTAAGAAT AATGTTGATT TCCCAAGTGC TTTACCTTAT 1603 CTGTTAAAGC GTAAGATGAA TTGGTATTTG CTTCATAGGC AGTTTGACTG 1651 CATGTATTAG AGAATGAAAA GAAGATATTT GTAGTAATGC CTGGAAACTT 1701 GGTGCTTTAA ATTAAGGTAC TCCTCTGCTG CTGTAGAATG GATTCCACAC 10 1751 AGTGGATAGC TATGGGTGAT TCAGAATATT ATGTTTAGAT TCCCATTTGT LBOL TAAGTTTATA AGTTTTGTGG GGAATTATGA ACTTACTGTG TACTACCTGC 1851 ATTTGTGCTG TGTGAAAAAT AAATACAAGG ATTCGTTTAG CTAATTCAAC 1901 TTACTACAAA GACAAATGTC TGTTTTTATT TGCCTGCTAG GATTGTCTTT 1951 TTTAAAAGTC ATTTTTATTT ATAGGAATAT GGGTGTTTCT ATAGGAAGAA 15 2003 ACAGGTTTTT TGTTTTTTGT TTTTTAAGAT AAATTTGACA AAGTTAACTG 2053 AAATTTATCT GGTCCATTTT ATTCATGCTA CTAAGATGGG AATCTTTAAA 2101 CACAAGGGTC AGCAAGCTTT GGCCCATGGA TTGGCCACCT GTTACGTAAA 2151 TAAAGTTTCT TTGAAACAAG CCTACACTCA TTCATTTATG TTTTGTCTGT 2201 GGTTGCTTTC CACAACTGCA GAGTTGTATG GCTTGCAAGT CTAAAAACAT 20 2251 TTACTATTTG GCCCTCTAAG AAAAAGTTAA GACACCTAGT CTAATGGCCT 2301 TTTGGGAAAA AACAAATCAC TAACTCATAA TCATTTATAT CCATTATTTT 2353 CTGCATAAAT GTAATGCTAT TGTACAGGGT TTGGTAGAAT AAATATTCAG 2401 ACTGACTAAA CTGTTCTAAA TTCTCACAAA AAAGTCCCCA AACAACATGC 2451 CTCCTAAAAA ACATTTTCCT ATCTTTTACA AGAGGTATGA ACATTTGTAG 2501 GGTTCCACAT TTGCATCTAG AAATCCAATG CTCTTTAGAA TGTTATTACG 25 2551 AATAGAAAGA TGGCCAGGAT GACCTTTAGT GTTACATGAT GTTCAGCAAA ZLOB TTTTAATTCA AACCTTGATA TGCCTGGACA CTGAAAAGTA AACGCATCAC 2653 CTCCTATTTT ATACCCTACC TTCTGGTTCC CAATTGGGAG AGCACATAGA 2701 GGGAAGGAGA CAATATAGAA ACTACGGAGT CCGCTGGTAG TGGGCTGCAT 30 2753 GGTGTGACAG AGCCCTTCTC TGTAAAATGG AAATGACACC ACTAGCCATC 2803 TCAATAGTTA CAAGAATTAA AAGAGATACA GTACCTGAAG TGCTTAGCGC 2851 ATGGTAGCAT TTCATAAATG TTTAGTGTCA ATACTAATGC TCTAATAATG 2903 TAAATTGTTA ATAATTTATT TCCCTAATAT CAGGAAATCC CAGTTGTCTA 295% TGTGGCCCAG TGCTTAAAAA CGCCTTCTTG CATGAGGGGA TTGAACTATA
300% CAATGTTTGT TAACTTTGTA TTTGTATTTT TTCCTATAAA ATCTTAAAAT 35 3051 AAAATTAGGA GATGTGTTCT GATGTAAAAA AAAAAAAAA AAAAAA

BLAST Results

40

Entry HS451B21 from database EMBL:
Human DNA sequence *** SEQUENCING IN PROGRESS *** from clone
451B21

45 Score = 11219, P = 0.0e+00, identities = 2287/2343

Medline entries

50

No Medline entry

55

Peptide information for frame 2

ORF from 20 bp to 940 bp; peptide length: 307 Category: similarity to unknown protein Classification: no clue

5 I MAALAPLPPL PAQLKSIQHH LRTAQEHDKR DPVVAYYCRL YAMQTGMKID
51 SKTPECRKFL SKLMDQLEAL KKQLGDNEAI TQEIVGCAHL ENYALKMFLY
101 ADNEDRAGRF HKNMIKSFYT ASLLIDVITV FGELTDENVK HRKYARWKAT
151 YIHNCLKNGE TPQAGPVGIE EDNDIEENED AGAASLPTQP TQPSSSSTYD
201 PSNMPSGNYT GIQIPPGAHA PANTPAEVPH STGVASNTIQ PTPQTIPAID
10 251 PALFNTISQG DVRLTPEDFA RAQKYCKYAG SALQYEDVST AVQNLQKALK
301 LLTTGRE

15 BLASTP hits

No BLASTP hits available

Alert BLASTP hits for DKFZphfbr2_78el8, frame 2

No Alert BLASTP hits found

Pedant information for DKFZphfbr2_78el8, frame 2

25

20

Report for DKFZphfbr2_78el8-2

ELENGTHI 313
30 EMWI 34463.75
EpII 5.64
EHOMOLI PIR:TO4778 hypothetical protein Flom23.70 Arabidopsis thaliana 3e-22
EKWI All_Alpha
35 EKWI LOW_COMPLEXITY -16.61 %

SEQ GEFGVEMAALAPLPPLPAQLKSIQHHLRTAQEHDKRDPVVAYYCRLYAMQTGMKIDSKTP SEG PRD **ECRKFLSKLMDQLEALKKQLGDNEAITQEIVGCAHLENYALKMFLYADNEDRAGRFHKNM** SEQ ------SEG PRD 45 SEQ **IKSFYTASLLIDVITVFGELTDENVKHRKYARWKATYIHNCLKNGETPQAGPVGIEEDND** SEG ----xxxxxx IEENEDAGAASLPTQPTQPSSSSTYDPSNMPSGNYTGIQIPPGAHAPANTPAEVPHSTGV 50 SEQ SEG PRD SEQ ASNTIQPTPQTIPAIDPALFNTISQGDVRLTPEDFARAQKYCKYAGSALQYEDVSTAVQN 55 SEG PRD SEQ LQKALKLLTTGRE

5 (No Prosite data available for DKFZphfbr2_7&el&.2)
(No Pfam data available for DKFZphfbr2_7&el&.2)

DKFZphfbr2_78i21

5 group: metabolism

DKFZphfbr2_7&i21 encodes a novel 477 amino acid protein with similarity to beta-aspartate methyltransferases.

- The L-isoaspartyl methyltransferase (Pimt), as an example, is a highly conserved enzyme utilising S-adenosylmethionine (AdoMet) to methylate aspartate residues of proteins damaged by agerelated isomerisation and deamidation.
- 15 The new protein can find application in diagnosis/modulation of protein damage and age-related degenerative processes.

unknown protein

20

weak similarity to beta-aspartate methyltransferase pimT of Mycobacterium leprae perhaps complete cds.

25 Sequenced by MediGenomix

Locus: unknown

Insert length: 1842 bp
30 Poly A stretch at pos- 1819, polyadenylation signal at pos- 1802

WO 01/98454 PCT/IB01/02050 LOCATTY TCAAGAGGAT GACCATGAAG AATCGCATTC TGATTTTCCA 1351 TATGGATCAT TTCCCTATGT TGCTAGACCA GTACACTGGC AACCTGGTCA 1401 TACAGCTTTT CTTGTCAAGT TGAGGAAGGT CAAACCACAA CTTAACTGAG 1451 TACTCCAGAT GACAGTAACT GACTTGAAGA TGGAAAAATA TCAAAATAGA 1501 ACTITATATI GAAAATCACI GCTTCCATAG ATTGGCATTT TTAGCTATTA 5 1551 CTATGACTTA TATAACTTAT ACATATAATT TTGAAAATAA CAACTAAAAG BEOB ATGTATAACA TAGCAAAACT GCTTAAACAT CCCATTTTGA CACTTGTCTT 1651 GCAGTTAGTT TGACATTTTG TAGTTAATGA TTCCAAATTG GTTTAGTTGG 1701 GCCATCTCAT TCTTCACTTC CTGTAAACCA CTCCATAGAT TTGTCTTTCT 1751 TCAAGAAATT AGTTTTCTTT CCTTTATTTG ATTGATGGTC ATTGACTACT 10 **BLAST Results** 15 No BLAST result 20 Medline entries No Medline entry 25 Peptide information for frame 1 30 ORF from 16 bp to 1446 bp; peptide length: 477 Category: putative protein Classification: no clue 1 MLMAWCRGPV LLCLRQGLGT NSFLHGLGQE PFEGARSLCC RSSPRDLRDG 35 51 EREHEAAQRK APGAESCPSL PLSISDIGTG CLSSLENLRL PTLREESSPR JOJ ELEDSSGDQG RCGPTHQGSE DPSMLSQAQS ATEVEERHVS PSCSTSRERP 151 FRAGELILAE TGEGETKFKK LFRLNNFGLL NSNWGAVPFG KIVGKFPGRI 203 LRSSFGKQYM LRRPALEDYV VLMKRGTAIT FPKDINMILS MMDINPGDTV 251 LEAGSGSGGM SLFLSKAVGS QGRVISFEVR KDHHDLAKKN YKHWRDSWKL 301 SHVEEWPDNV DFIHKDISGA TEDIKSLTFD AVALDMLNPH VTLPVFYPHL 40 351 KHGGVCAVYV VNITQVIELL DGIRTCELAL SCEKISEVIV RDWLVCLAKQ 401 KNGILAQKVE SKINTDVQLD SQEKIGVKGE LFQEDDHEES HSDFPYGSFP 451 YVARPVHWQP GHTAFLVKLR KVKPQLN 45 BLASTP hits

No BLASTP hits available

50

Alert BLASTP hits for DKFZphfbr2_78i21, frame 1

No Alert BLASTP hits found

Pedant information for DKFZphfbr2_78i2l, frame l

Report for DKFZphfbr2_78i21.1

```
ELENGTHD
         482
          53521.20
   EWWI
5
   [[q]
          6.58
             TREMBL: AFD88800_2 product: "unknown"; Rhodococcus
   [HOMOL]
   erythropolis ARC (arc) gene, complete cds; and unknown genes. 2e-
   EFUNCATD r general function prediction
                                 IM. jannaschii:
10
   MJ01341 6e-10
   EFUNCATE 05.07 translational control
                              CS. cerevisiae, YJL125c1
   6e-04
   EBFOCKZI
         BLOOSOJE
   EBFOCK23
         BL01279A
15
          Alpha_Beta
   CKWI.
   EKWI
          LOW_COMPLEXITY
                       2.49 %
   SEQ
      PSRNTMLMAWCRGPVLLCLRQGLGTNSFLHGLGQEPFEGARSLCCRSSPRDLRDGEREHE
20
   SEG
   PRD
      SEQ
      AAQRKAPGAESCPSLPLSISDIGTGCLSSLENLRLPTLREESSPRELEDSSGDQGRCGPT
   SEG
25
   PRD
      SEQ
      HQGSEDPSMLSQAQSATEVEERHVSPSCSTSRERPFQAGELILAETGEGETKFKKLFRLN
   SEG
   PRD
      30
   SEQ
      NFGLLNSNWGAVPFGKIVGKFPG@ILRSSFGK@YMLRRPALEDYVVLMKRGTAITFPKDI
   SEG
   PRD
      35
      NMILSMMDINPGDTVLEAGSGSGGMSLFLSKAVGSQGRVISFEVRKDHHDLAKKNYKHWR
   SEQ
   SEG
      PRD
      SEQ
      DSWKLZHVEEWPDNVDFIHKDISGATEDIKSLTFDAVALDMLNPHVTLPVFYPHLKHGGV
40
  SEG
  PRD
      SEQ
      CAVYVVNITQVIELLDGIRTCELALSCEKISEVIVRDWLVCLAKQKNGILAQKVESKINT
  SEG
45
  PRD
      SEQ
      DVQLDSQEKIGVKGELFQEDDHEESHSDFPYGSFPYVARPVHWQPGHTAFLVKLRKVKPQ
  SEG
  PRD
      50
  SEQ
      LN
  SEG
  PRD
      CC
55
   (No Prosite data available for DKFZphfbr2_78i21.1)
   (No Pfam data available for DKFZphfbr2_78i21.1)
```

DKFZphmel2_12jl

10

20

5 group: melanoma derived

DKFZphmel2_12jl encodes a novel 905 amino acid protein, which has similarity to integrin I of Saccharomyces cerevisiae.

The novel protein contains a leucin zipper.
No informative BLAST results: No predictive prosite, pfam or SCOP motife.

15 The new protein can find application in studying the expression profile of melanoma-specific genes.

weak similarity to integrin I (Saccharomyces cerevisiae)

Sequenced by EMBL

Locus: unknown

25 Insert length: 2942 bp
Poly A stretch at pos. 2926, no polyadenylation signal found

1 CGAAAGCTAA AGGCCGGCGC ACGCTGGGCG GTGGTGGTCC CTAAGCCGGG 30 51 CCGCGGCCGG TGCAATGGAC TCCACTGCCT GCTTGAAGTC CTTGCTCCTG 101 ACTGTCAGTC AGTACAAAGC CGTGAAGTCA GAGGCGAACG CCACTCAGCT 151 TTTGCGGCAC TTGGAGGTAA TTTCTGGACA GAAACTCACA CGACTATTTA 201 CATCAAATCA GATATTAACA AGTGAATGCT TGAGTTGCCT TGTAGAGCTA 251 CTTGAAGACC CCAACATAAG TGCTTCACTG ATCTTAAGTA TTATCGGTTT BOL GCTGTCTCAA CTAGCAGTAG ACATTGAAAC CAGAGATTGT CTTCAGAATA 35 351 CATATAATCT GAATAGTGTG CTGGCGGGAG TGGTTTGTCG GAGCAGCCAC 401 ACTGATTCGG TGTTTTTGCA GTGCATTCAA CTTCTACAGA AGTTAACATA 451 TAATGTCAAA ATTTTCTATT CTGGTGCCAA TATAGATGAA TTAATTACGT 501 TCCTGATAGA TCACATTCAA TCTTCTGAAG ATGAGTTAAA AATGCCTTGT 551 CTAGGATTAT TGGCAAATCT TTGTCGGCAC AATCTTTCTG TTCAAACGCA 40 LOD CATAAAGACA TIGAGTAATG TGAAATCTTT TTATCGAACT CTTATCACCT **L51 TGTTGGCCCA TAGTAGTTTA ACTGTGGTTG TGTTTGCACT TTCAATATTA** 701 TCCAGTTTGA CATTAAATGA AGAGGTGGGG GAAAAGCTAT TCCATGCTCG 751 AAACATTCAT CAGACTTTTC AACTAATATT TAATATTCTC ATAAACGGTG 45 BOD ATGGCACTCT AACTAGAAAG TATTCAGTTG ACCTACTGAT GGATCTCCTT BS% AAGAATCCTA AAATTGCTGA TTATCTCACC AGATATGAGC ACTTTTCTTC 9D1 ATGTCTTCAC CAAGTATTAG GTCTTCTTAA TGGAAAGGAT CCTGATTCCT 951 CTTCAAAGGT TTTAGAATTA CTTCTTGCCT TCTGTTCAGT GACTCAGCTG DODD CGCCATATGC TCACTCAGAT GATGTTTGAA CAGTCTCCAC CTGGCAGCGC 50 3053 CACTCTGGGA AGCCATACTA AATGTTTAGA ACCTACTGTG GCTCTACTGC BIOI GCTGGTTAAG CCAACCTTTG GACGGATCAG AAAACTGTTC TGTTTTAGCA 1151 TTGGAGTTGT TCAAGGAAAT ATTTGAGGAT GTCATAGATG CTGCTAACTG 1201 TTCCTCGGCT GATCGTTTTG TGACCCTTCT GCTGCCTACA ATCCTTGATC 1251 AACTTCAGTT CACAGAACAA AATCTAGATG AGGCTTTAAC AAGAAAAAT 55 1301 GTGAAAGGGA TTGCCAAGGC CATTGAAGTT TTGTTAACTC TCTGTGGAGA 1351 TGATACACTA AAAATGCATA TTGCAAAAAT CTTGACAACT GTCAAGTGTA 1401 CCACTCTTAT AGAACAACAA TTTACATATG GCAAGATTGA CCTGGGATTT 1451 GGAACAAAGG TTGCAGATTC TGAATTATGC AAACTTGCTG CTGATGTAAT

	•		
j.			
		•	
			,

1501 TTTGAAAACT CTTGATTTGA TTAACAAACT TAAACCATTG GTTCCTGGTA
1551 TGGAAGTAAG CTTCTACAAA ATACTTCAGG ACCCACGTTT GATTACTCCT
1601 TTGGCTTTTG CTTTAACGTC AGATAATAGA GAACAAGTAC AGTCTGGACT
1651 GAGAATATTA TTGGAGGCTG CTCCACTGCC AGATTTTCCT GCTTTAGTAC
1701 TTGGAGAAAG TATAGCAGCA AACAATGCCT ATAGACAACA GGAAACAGAA
1751 CATATACCCA GAAAAATGCC CTGGCAATCA TCAAATCACA GTTTTCCAAC
1801 ATCAATAAAG TGTTTAACTC CTCATTTGAA AGATGGTGTT CCTGGATTGA 5 LB5L ATATTGAAGA ATTAATAGAG AAACTTCAGT CTGGAATGGT GGTAAAGGAT LADI CAGATTIGIG ATGIGAGAAT ATCIGACATA ATGGATGTAT ATGAAATGAA 1951 ACTATCCACA TTAGCTTCCA AAGAAAGCAG GCTACAAGAT CTTTTGGAAA 10 2001 CAAAAGCTCT AGCCCTTGCA CAGGCTGATA GACTGATTGC TCAGCATCGC 2051 TGTCAAAGAA CTCAAGCTGA AACAGAGGCA CGGACACTTG CTAGTATGTT 2101 GAGAGAGTT GAGAGAAAAA ATGAAGAGCT TAGTGTGTTG CTGAAGGCGC 2151 AGCAAGTTGA ATCAGAAAGA GCGCAGAGTG ATATTGAGCA TCTCTTTCAA 2201 CATAATAGGA AGTTAGAGTC TGTGGCTGAA GAACATGAAA TACTGACAAA 15 225% ATCCTACATG GAACTTCTTC AGAGAAATGA AAGTACTGAA AAGAAGAATA DTDADADAD ADATTADA ADTOTOTAD TOTADADADA ADES 2351 AAAAAATTA ATGAGTCACT CAAGGAACAA AATGAAAAA GTATGCCCA 2401 ATTAATAGAG AAAGAAGAAC AGAGAAAAGA AGTACAGAAT CAGCTAGTAG 245% ACAGAGAACA TAAGCTAGCA AATTTGCATC AAAAAACAAA AGTACAAGAA 20 2501 GAAAAGATTA AAACCTTACA AAAGGAAAGG GAAGATAAGG AAGAAACCAT 2551 TGATATCCTT AGAAAAGAAT TAAGCAGAAC AGAACAGATA AGAAAAGAGT ZLOL TGAGCATTAA GGCTTCCTCC CTAGAGGTTC AAAAGGCACA ATTAGAAGGT 2651 CGTTTGGAAG AGAAAGAGTC CTTGGTGAAA CTTCAGCAAG AGGAATTGAA 2701 CAAACACTCC CACATGATAG CAATGATCCA CAGTTTAAGT GGTGGAAAAA 25 2751 TAAATCCAGA AACTGTGAAT CTCAGTATAT AGACATTATG GCATTTTGGA 2801 ATTTGTAATC TCATGATATT TTTGATGTAT TTATCTATTG GAGGGGGGGT 2851 GGGTAGGGGA GTTAATTTGT GACTTCGTAA CAATAAGAAG TTATTATCTA 2901 ATTTAGTAAA GACCCTGATC TGTTGCAAAA AAAAAAAAA AA 30

BLAST Results

35 No BLAST result

Medline entries

40 95039111:

Hostetter MK, Tao NJ, Gale C, Herman DJ, McClellan M, Sharp RL, Kendrick KE, Antigenic and functional conservation of an integrin

45 I-domain in

55

Saccharomyces cerevisiae. Biochem Mol Med 1995 Aug;55(2):122-30

19458454:

Berton G. Lowell CA.; Integrin signalling in neutrophils and macrophages. Cell Signal 1999 Sep;11(9):621-35

Peptide information for frame 2

ORF from 65 bp to 2779 bp; peptide length: 905

Category: putative protein

Classification: Cellular transport and traffic

Prosite motifs: LEUCINE_ZIPPER (331-352)

5 1 MDSTACLKSL LLTVSQYKAV KSEANATQLL RHLEVISGQK LTRLFTSNQI 51 LTSECLSCLV ELLEDPNISA SLILSIIGLL SQLAVDIETR DCLQNTYNLN 101 SVLAGVVCRS SHTDSVFLQC IQLLQKLTYN VKIFYSGANI DELITFLIDH 151 IQSSEDELKM PCLGLLANLC RHNLSVQTHI KTLSNVKSFY RTLITLLAHS

201 SLTVVVFALS ILSSLTLNEE VGEKLFHARN IHQTFQLIFN ILINGDGTLT 10 251 RKYSVDLLMD LLKNPKIADY LTRYEHFSSC LHQVLGLLNG KDPDSSSKVL

301 ELLLAFCSVT QLRHMLTQMM FEQSPPGSAT LGSHTKCLEP TVALLRWLSQ 351 PLDGSENCSV LALELFKEIF EDVIDAANCS SADRFVTLLL PTILDQLQFT 401 EQNLDEALTR KNVKGIAKAI EVLLTLCGDD TLKMHIAKIL TVKCTTLIE

451 QQFTYGKIDL GFGTKVADSE LCKLAADVIL KTLDLINKLK PLVPGMEVSF 501 YKILQDPRLI TPLAFALTSD NREQVQSGLR ILLEAAPLPD FPALVLGESI 551 AANNAYRQQE TEHIPRKMPW QSSNHSFPTS IKCLTPHLKD GVPGLNIEEL 601 IEKLQSGMVV KDQICDVRIS DIMDVYEMKL STLASKESRL QDLLETKALA

651 LAQADRLIAQ HRCQRTQAET EARTLASMLR EVERKNEELS VLLKAQQVES

701 ERAGSDIEHL FRHNRKLESV AEEHEILTKS YMELLARNES TEKKNKDLRI 751 TCDSLNKRIE TVKKLNESLK ERNEKSIARL IEKEERRKEV RNALVDREHK 801 LANLHRKTKV REEKIKTLRK EREDKEETID ILRKELSRTE RIRKELSIKA 20 851 SSLEV@KA@L EGRLEEKESL VKL@@EELNK HSHMIAMIHS LSGGKINPET

POT ANTZI

15

25

45

50

BLASTP hits

30 No BLASTP hits available

Alert BLASTP hits for DKFZphmel2_l2jl, frame 2

TREMBL:SCINTANA_1 Saccharomyces cerevisiae integrin analoque 35 gene, complete cds., N = 1, Score = 216, P = 1.3e-13

>TREMBL:SCINTANA_1 Saccharomyces cerevisiae integrin analogue 40 gene, complete cds.

Length = $1_{1}015$

HSPs:

Score = 216 (32.4 bits), Expect = 1.3e-13, P = 1.3e-13 Identities = 80/302 (26%), Positives = 155/302 (51%)

597 IEELIEKLQSGMVVKDQICDVRISDIM---EE3 DAJAJAJJUDJRZBAZKALTKALALAQ 653

> I L EKL++ D+ ++2I++ + E +L+

AL +

275 ISLLKEKLETATTANDENVN-Sbict:

KISELTKTREELEAELAAYKNLKNELETKLETSEKALKE 333

55 **L54 A---DRLIAQHRCQRTQAETEAR----TLASMLREVERKNEELSVLLKA--**Query: QQVESERAQ 704

```
+ Q + TE +
                                            +L + L +E+++E+L+ LK
    +Q+ ++
    Sbict:
            334
    VKENEEHLKEEKIQLEKEATETKQQLNSLRANLESLEKEHEDLAAQLKKYEEQIANKERQ 393
 5
             7D5 SDIEHLFQHNRKLESVAEEHEILTKSYMEL---LQRNESTEKKNKDLQIT-
    Querv:
    CDSLNKQIE 760
                  + E + Q N ++ S +E+E + K
                                             EL
                                                     + + + + TZ+
    D+LN QI+
    Sbjct:
             -NY PPE
10
    EEISQLNDEITSTQQENESIKKKNDELEGEVKAMKSTSEEQSNLKKSEIDALNLQIK 452
    Query:
             763
    TVKKLNESLKEQNEKSIAQLIEKEEQRKEVQNQLVDREHKLANLHQKTKVQEEKIKT--- 817
                               +SI + + + KE+Q++
15
                  +KK NE+ +
                                                     +E +++ L K K
    E+K
    Sbjct:
             453
    ELKKKNETNEASLLESIKSIESETVKIKELQDECNFKEKEVSELEDKLKASEDKNSKYLE 575
20
    Query:
             818 LQKEREDKEETIDI----LRKELSRTEQIRKELSIKASSLE-
    VQKAQLEGRLEEKESLVK 872
                 LQKE E +E +D
                                  L+"+L +
                                             + K
                                                      S L ++K
                                                               ER
    +E L K
    Sbjct:
             513
25
    LQKESEKIKEELDAKTTELKIQLEKVTNLSKAKEKSESELSRLKKTSSEERKNAEEQLEK 572
    Query:
             873 LQQE 876
                 L+ E
             573 LKNE 576
    Sbjct:
30
     Score = 186 (27.9 bits), Expect = 2.0e-10, P = 2.0e-10
     Identities = 82/301 (27%), Positives = 155/301 (51%)
             598 EELIEKLQSGMVVKDQICDVRISDIMDVYEMKLSTLASKESR---LQD-
35
    LLETKALALAQ 653
                 +ELI +LQ+
                             +K +
                                       ++2
                                                      K++
                                                              LQD +L
    Sbjct:
             PII DELI-
    RLQNENELKAKEIDNTRSELEKVSLSNDELLEEKQNTIKSLQDEILSYKDKITRN 669
40
    Query:
    ADRLIAGHRCGRTGAETEARTLASMLREVERKNEELSVLLKAGGVESERAGSDIEHLFGH 713
                  ++L++ R +
                                E+
                                      L
                                         LR +
                                                        LK + ES +
45
    Sbict:
            670 DEKLLSIERDSKRDLES----
    LKEGLRAAGESKAKVEEGLKKLEEESSKEKAELEKSKEM 725
            714 NRKLESVAEEHEILTKSYMELLQRN-ESTEKKNKDLQITCDSL-
    Query:
   NKQIETVKKLNESLKE 771
50
                  +KLEZ E +E
                                KS ME ++++ E E+ K +
    ++NES K+
    Sbjct:
            726
    MKKLESTIESNETELKSSMETIRKSDEKLEQSKKSAEEDIKNLQHEKSDLISRINESEKD 785
            772 QNE-KSIAQLIEKEEQRKE-VQNQLVDREHKL-
55
    Query:
    ANLHQKTKVQEEKIKTLQKEREDKEET 828
                                   E V+ +L + + K+
                  E KZ
                        ++ K
                                                   Ν
                                                      + T V + K++
   +++E +DK+
```

Sbjct: 786 IEELKSKLRIEAKSSSELETVKQELNNAQEKIRVNAEENT-VLKSKLEDIERELKDKQAE 844

Query: B29 IDILR--KEL--SRTEQIRKEL-----SIKASSLEVQKAQLE5 GRLEEKESLVKLQ B74

I + KEL SR +++ +EL

Z + Z EV+K Q+E

Sbjct: 845

IKSNQEEKELLTSRLKELEGELDSTQQKAQKSEEESRAEVRKFQVEKSQLDEKAMLLETK 904

10

Query: 875 QEEL-NK 880

+L NK

Sbjct: 905 YNDLVNK 911

15 Score = 173 (26.0 bits), Expect = 5.7e-09, P = 5.7e-09 Identities = ??/287 (26%), Positives = 146/287 (50%)

Query: 601 IEKLQSGMVVKDQICDVRISDIMDVYEMKLSTLASKES--RLQDLLETKALALAQADRLI 658

20 ++K + + K++ + IS + D E+ ST ES + D LE +
A+
Sbjct: 38D LKKYEEQIANKERQYNEEISQLND--EITSTQQENESIKKKNDELEGEVKAMKST---- 432

25 Query: LS9 AQHRCQRTQAETEARTLASMLREVERKNE-ELSVLLKAQQVESERAQSDIEHLFQH-NR 715

++ + ++E +A L ++E+++KNE E S+L + +ESE + I+

_ · N

Sbjct: 433 SEEQSNLKKSEIDALNL--QIKELKKKNETNEASLLESIKSIESETVK--IKELQDECNF 488

Query: 716 KLESVAEEHEILTKSY--MELLQRNESTEKKNKDLQITCDSLNKQIETVKKLNESLKEQ 772

K + V+E + L Z + L+ + + EK ++L L Q+E V

35 L+++ KE+ Sbjct: 489

KEKEVSELEDKLKASEDKNSKYLELQKESEKIKEELDAKTTELKIQLEKVTNLSKA-KEK 547

Query: 773 NEKSIAQLIE-KEEQRKEVQNQL--VDREHKLAN--

40 LHQKTKVQEEKIKTLQKEREDKEE 827

+E +++L + E+RK + QL + E ++ N ++ K+ E T+

+E +K

Sbjct: 548

SESELSRLKKTSSEERKNAEE@LEKLKNEI@IKN@AFEKERKLLNEGSSTIT@EYSEKIN LO?

45

30

Query: 828 TI-

DILRKELSRTEQIRKELSIKASSLEVQKAQLEGRLEEKESLVKLQQEELNKHSHMI 885
T+ D L + + E KE+ S LE + LEEK++ +K Q+E+

+ I

50 Sbjct: 608

TLEDELIRLQNENELKAKEIDNTRSELEKVSLSNDELLEEKQNTIKSLQDEILSYKDKI LLL

Score = 171 (25.7 bits), Expect = 9.3e-09, P = 9.3e-09 Identities = 76/311 (24%), Positives = 152/311 (48%)

55

Query: 596 NIEELIEKLQSGMVVKDQ------ICDVRISDIMDVYEMKLSTLASKESRLQDLLETKA 648

WO 01/98454 PCT/IB01/02050 N EE +EKL++ + +K+QI 2 + +K++TL + RLQ+ E KA Sbjct: 565 NAEEGLEKLKNEIGIKNGAFEKERKLLNEGSSTITGEYSEKINTLEDELIRLQNENELKA 624 5 649 LALAQADRLIAQHRCQRTQA-ETEARTLASMLREVERKNEELSVL-LKAQQVESERAQSD 706 + E + T + S + E +K +E + ++ D Sbjct: 625 10 KEIDNTRSELEKVSLSNDELLEEKQNTIKSLQDEILSYKDKITRNDEKLLSIERD-SKRD 683 707 IEHLFQHNRKL-ESVAEEHEILTKSYMELLQRNESTEKKN---KDLQITCDS----LNKQ 758 15 +EL+RES A+ E L K Ε EK K.L+T+SSbjct: 684 LESTKEGT BY BY AND THE STATE OF 20 Query: 759 IETVKKLNESLKEQNEKSIAQLIEK-EERKEVRNALVDREHKLANLHRKTKVREE---K 814 +ET++K +E L EQ++KS + I+ + ++ ++ +++ + E Sbict: 744 METIRKSDEKL-25 EGZKKZAEEDIKNLGHEKZDLIZRINEZEKDIEELKZKLRIEAKZZZE 805 815 IKTLQKEREDKEETIDILRKE----LSRTEQIRKELSIKASSL---EVQKAQLEGRLEEK 867 ++T+++E + +E I + +E S+ E I +EL K + ++ +K 30 L RL+E Sbjct: 803 LETVKQELNNAQEKIRVNAEENTVLKSKLEDIERELKDKQAEIKSNQEEKELLTSRLKEL ALZ 8P8 EZFAKTőGEEFNK 990 Query: 35 + Q++ K. Sbjct: 863 EQELDSTQQKAQK 875 Score = 165 (24.8 bits), Expect = 4.1e-08, P = 4.1e-08 Identities = 65/286 (22%), Positives = 149/286 (52%) 40 595 LNIEELIEKLQSGMVVKDQICDVR-ISDIMDVYEMKLSTLASKESRL-Query: QDLLETKALALA 652 +N ++ + L+ + K I +++ I++ ++ +++ + L+ ++ ++L+E K+ 114 VNHQKETKSLKEDIAAK--45 Sbict: ITEIKAINENLEKMKIQCNNLSKEKEHISKELVEYKS-RFQ 170 653 QADRLIAQHRCQRTQAETEARTLASMLREVERKNEELSVLLKAQQVESE---Query: -RAQSDIE 708 50 D L+A+ T+ + ++LA+ ++++ +NE L + EZ Q+ I+ Sbict: 171 SHDNLVAK----LTE---KTKZTVNJKDMGVENEZTIKVAEEZKNEZZIGTZNTGVKID 553 55 709 HLFQH--NRKLE--Query: SVAEEHEILTKSYMELLQRNESTEKKNKDLQITCDSLNKQIETVKK 764

N ++E S++ELK++LQE

K+

D

+ Q

QT

+K+

Sbjct: 224 SMSQEKENFQIERGSIEKNIEQLKKTISDLEQTKEEIISKSDSSK--DEYESQISLLKE 280

Query: 765

Sbjct: 281

KLETATTANDENVNKISELTKTREELEAELAAYKNLKNELETKLETSEKALKEVKENEEH 340

10

Query: 825 KEETIDILRKELSRTEQIRKELSIKASSLEVQKAQLEGRLEEKESLVKLQQEELNK 880
KEE I L KE + T+Q L SLE + L +L++ E + ++
+ N+

15 Sbjct: 341 LKEEKIQ-

LEKEATETKQQLNSLRANLESLEKEHEDLAAQLKKYEEQIANKERQYNE 396

Score = 158 (23.7 bits), Expect = 1.9e-07, P = 1.9e-07 Identities = 74/268 (27%), Positives = 136/268 (50%)

20

30

25 Sbjct: 695

@ESKAKVEEGLKKLEEESSKEKAELEKSKEMMKKLESTIESNETELKSSMETIRKSDEKL 754

Query: L53
QADRLIAGHRCGRTGAETEARTLASMLREVERKNEELSVLLKAGGVESERAGSDIEHLFG 732
+ + A+ + Q E L S + E E+ EEL L+ + S

+ + A+ + Q E L S + E E-++ + L Sbjct: 755 EQSKKSAEEDIKNLQHEKS--DLISRINESEKDIEELKSKLRIEAKSSSELETVKQELNN 812

35 Query: 713 HNRKLESVAEEHEILTKSYMELLQRNESTEKKNKDLQITCDSLNKQIET--VKKLNESLK 770

K +

+K+L + L Sbjct: 813 AQEKIRVNAEENTVL-KSKLEDIER----

40 ELKDKQAEIKSNQEEKELLTSRLKELEGELD 867

Query: 771 EQNEKSIAQLIEKEEQRKEVQNQLVDR---EHKLANLHQKTKVQEEKIKTLQKEREDKEE 827 +K AQ E EE R EV+ V++ +

+K AQ E EE R EV+ V++ + K L K

45 +++ + ++
Sbjct: 868 STQQK--AQKSEEESRAEVRKFQVEKSQLDEKAMLLETKYNDLVNKEQAWKRDEDTVKK 924

Query: 828 TIDILRKELSRTEQIRKEL-SIKASSLEVQKAQLEGRLE 865 T D R+E+ E++ KEL ++KA + ++++A E R E Sbjct: 925 TTDSQRQEI---EKLAKELDNLKAENSKLKEAN-EDRSE 959

Score = 155 (23.3 bits), Expect = 3.9e-07, P = 3.9e-07

Score = 155 (23.3 bits), Expect = 3.9e-07, P = 3.9e-07 Identities = 73/269 (27%), Positives = 133/269 (49%)

55

Query: 624 DVYEMKLSTLASKESRLQD-LLETKALALAQADRLIAQHRCQRTQAET--EARTLASML 679

```
LQD +L K
                                               ++L++ R + E+
                ++ E K +T+ S
    R
    Sbjct:
           PA3 EFFEEKGNLIKZ----
    LQDEILSYKDKITRNDEKLLSIERDSKRDLESLKEQLRAAQESK 698
    Query: 680 REVE---
    RKNEELSVLLKAQQVESERAQSDIEHLFQHNRKLESVAEEHEILTKSYMELLQ 736
                +VE +K EE 2 KA+ +S+
                                             +E
                                                  + N
10
    Sbjct: b99 AKVEEGLKKLEEESSKEKAELEKSKEMMKKLESTIESNET--
    ELKSSMETIRKSDEKLE@ 756
    Query: 737 RNESTEKKNKDLQITCDSLNKQIETVKKLNESLKEQ---
    NEKSIAQLIEKEEQRKEVQNQ 793
15
                  +S E+ K+LQ
                                  L +I +K E LK +
                                                        KZ ++L
    Sbjct:
            757
    ZKKZAEEDIKNL@HEKZDLIZRINESEKDIEELKZKLRIEAKZZZELETVK@ELNNA@EK 47P
20
    Querv: 794 L-VDREH-----
    KLANLHQKTKVQEEKIKTLQKEREDKEETIDILRKELSRTEQIRKEL 846
                + V+ E
                            KL ++ ++ K ++ +IK+ Q+E+E
    T+Q + +
    Sbjct:
            817
25
   IRVNAEENTVLKSKLEDIERELKDKQAEIKSNQEEKELLTSRLKELEQELDSTQQ-KAQK 875
            847 SIKASSLEVQKAQLE-GRLEEKESLVKLQQEEL-NK 880
    Query:
                S + S EV+K Q+E +L+EK L++ + +L NK
            876 SEEESRAEVRKFQVEKSQLDEKAMLLETKYNDLVNK 911
    Sbjct:
30
     Score = 146 (21.9 bits), Expect = 3.5e-06, P = 3.5e-06
     Identities = 73/311 (23%), Positives = 152/311 (48%)
    Query: 520 DNREQVQSGLRIL----LEAAPLPDFPALV--
35
   LGESIAANNAYRQQETEHIPRK-MPWQ 571
                +++ +V+ GL+ L
                                  EAL
                                            ++ L
                                                  +I +N
                                                                EI
   Sbict:
            696
   ESKAKVEEGLKKLEEESSKEKAELEKSKEMMKKLESTIESNETELKSSMETIRKSDEKLE 755
40
   Query: 572 SSNHSFPTSIKCLTPHLKDGVPGLNIEEL-
   IEKLQSGMVVKDQICDVRISDIMDVYEMKL 630
                         IK L
                                 D + +N E IE+L+S + +
    ++ + +L
45
   Sbjct: 75b QSKKSAEEDIKNLQHEKSDLISRINESEKDIEELKSKLRI-----
   EAKZZZELETVKŒEL 810
   Query: L31 STLASK---
   ESRLQDLLETKALALAQADRLIAQHRCQRTQAETEARTLASMLREVERKNE 687
50
                    K
                         +
                               +L++K
                                      L +R +
                                                + +
                                                        + E
                                                              L Z
   L+E+E++ +
   Sbjct: Bll NNAQEKIRVNAEENTVLKSK---
   LEDIERELKDKQAEIKSNQEEKELLTSRLKELEQELD 867
   Query: 688
   ELSVLLKAQQVESERAQSDIEHLFQHNRKLESVAEEHEILTKSYMELLQRNESTEKKNKD 747
                     KAQ+ E E +++++ FQ + + E+ +L Y +L+ +
```

Sbjct: 868 --STQQKAQKSEEE-SRAEVRK-FQVEKS--QLDEKAMLLETKYNDLVNKEQAWKRDEDT 923

duery: 748 LQITCDSLNKQIETVKKLNESLKEQNEKSIAQLIEKEEQRKEVQNQLV--5 DREHKLANL 804

++ T DZ ++IE + K ++LK +N K LE E R E+ + ++ D

Sbjct: 922 VKKTTDSQRQEIEKLAKELDNLKAENSK----LKEANEDRSEIDDLMLLVTDLDEK--NA 975

10

Query: 805 HQKTKVQEEKIKTLQKEREDKEETID 830
++K+++ ++ E +D+EE D

Sbjct: 976 KYRSKLKDLGVEISSDEEDDEEEEDD 1001

15 Score = 146 (21.9 bits), Expect = 4.6e-06, P = 4.6e-06 Identities = 82/313 (26%), Positives = 145/313 (46%)

Query: 598

EELIEKL@SGMVVKD@ICDVRISDIMDVYEMKLSTLASKESRL@DLLETKALALA@ADRL 657
20 EEL +L + +K+++ + E+K + KE ++@ LE +A

Sbjct: 304 EELEAELAAYKNLKNELETKLETSEKALKEVKENEEHLKEEKIQ--LEKEATETKQQ--- 358

25 Query: 658 IAQHRCQRTQAETEARTLASMLREVERK-----NEELSVL--LKAQQVESERAQSD 706

+ R E E LA+ L++ E + NEE+S L + + Q

Sbjct: 359

E+E +

30 LNSLRANLESLEKEHEDLAAQLKKYEEQIANKERQYNEEISQLNDEITSTQQENESIKKK 418

Query: 707 IEHLFQHNRKLESVAEEHEILTKSYMELLQRN-ESTEKKNKDLQITCDSLNKQIET-VKK 764

+ L + ++5 +EE L K2 ++ L + +KKN+ + + K

E +

+L

K+

L KT

35 IE+ K
Sbjct: 419
NDELEGEVKAMKSTSEEQSNLKKSEIDALNLQIKELKKKNETNEASLLESIKSIESETVK 47B

Query: 765 LNESLKEQN--EKSIAQLIEK---EEQRKEVQNQLVDREHKLAN-LHQKT--40 -KVQEEKI 815

-KAMFFKI 972

K+Q EK+ Sbjct: 479

IKELQDECNFKEKEVSELEDKLKASEDKNSKYLELQKESEKIKEELDAKTTELKIQLEKV 538

45

Query: 816

KTLQKEREDKEETIDILRKELSRTEQIRKELSIKASSLEVQKAQLEGRLEEKESLVKLQQ 875 L K +E E ELSR ++K S + + E Q +L+ ++ K

+ ++

50 Sbjct: 539 TNLSKAKEKSES-----ELSR--LKKTSSERKNAEEQLEKLKNEIQIKNQAFEKER 588

+ E

Query: 876 EELNKHSHMIAMIHSLSGGKINPETVNL 903

+ LN+ Z I +Z + E + L

E N EK +++L +K

55 Sbjct: 589 KLLNEGSSTITQEYSEKINTLEDELIRL 616

Score = 145 (21.8 bits), Expect = 5.9e-06, P = 5.9e-06 Identities = 59/246 (23%), Positives = 115/246 (46%)

Query: 634 ASKESRLQ-DLLETKALALAQADRLIAQHRCQRTQAETEARTLASMLREVERKNEELSVL 692 + ES +Q L+ K +++Q + +R E + ++E+ Sbjct: 503 ZKNEZZIĞPZNPĞNKIDZWZĞEKE---NFQIERGSIEKNIEQLKKTISDLEQTKEE--II 261 693 LKAQQVESERAQSDIEHLFQHNRKLESVAEEHEI----10 LTKSYMELLQRNESTEKKNKD 747 LTK+ EL K+ + E +S I L + + + A + + Sbict: 2P5 ZKZDZZKDEA-EZ6IZ-LLKEKLETATTANDENVNKISELTKTREELEAELAAYKNLKNE 319 15 Query: 748 LQITCDSLNKQIETVKKLNESLKEQNEKSIAQLIEKEEQRKEVQNQLVDREHKLANLHQK 807 ++ K ++ VK+ E LKE+ + + E ++Q ++ L L+ 20 Sbjct: 320 LETKLETSEKALKEVKENEEHLKEEKIQLEKEATETKQQLNSLRANLESLEKEHEDLAAQ 379 Query: 808 TKVQEEKIKTLQKEREDKEETIDILRKELSRTEQIRKELSIKASSLEVQKAQLEGRLEEK 867 25 K EE+I KER+ EE I L E++ T+Q + + K LE + 380 LKKYEEQIAN--KERQYNEE-Sbict: ISQLNDEITSTQQENESIKKKNDELEGEVKAMKSTSEEQ 436 30 Query: 868 ESLVKLQQEELN 879 +L K + + LNSbjct: 437 SNLKKSEIDALN 448 Score = 137 (20.6 bits), Expect = 4.2e-05, P = 4.2e-05 35 Identities = 81/312 (25%), Positives = 140/312 (44%) 598 EELIEKLQSGMVVKDQICDVRISDIMDVYEMKLSTLASK-ESRLQDLLET-KALALAQAD 655 +EL ++++ ++ +++ Z+I D +++ L K E+ 40 K++ 420 DELEGEVKAMKSTSEEQSNLKKSEI-DALNLQIKELKKKNETNEASLLESIKSIESETVK 478 Query: 656 45 RLIAGHRCORTGAETEARTLASMLREVERKNEELSVLLKAQQVESERAQSDIEHLFQHNR 715 a c EE L L+ EKN + LK + E Sbict: 479 IKELQDECNFK--EKENZETEDKTKVZEDKNZKATETŐKEZEKIKEETDVKLLETKIŐTE 23P 50 Query: 716 KLESVAEEHEILTKSYMELLQRNESTEKKNKDLQITCDSLNKQIETVKKLNESLKEQNEK 775 K+ ++++ E ++S + L++ S E+KN + Q+ QI+ + 55 Sbjct: 537 KVTNLSKAKE-KSESELSRLKKTSSEERKNAEEQLEKLKNEIQIKN-QAFEKERKLLNEG 594

·776 SIAQLIEKEEQRKEVQNQLV--DREHKL-ANLHQKTKVQEEKIKTLQKER-Querv: EDKEETIDI 831

++++L+ E++L A E E+ T+ + EK+

E+K+ TI

Sbjct: 595

Query: 832 LRKE-LSRTEQI----RKELSIKASS---LEVQKAQLEGRLEEK----ESLVKLQQE--- 876

10 L+ E LS ++I K rzi+ z LE K QL ΕK E L

KL++E

Sbjct: 655

LQDEILSYKDKITRNDEKLLSIERDSKRDLESLKEQLRAAQESKAKVEEGLKKLEEESSK 714

877 ---ELNKHSHMIAMIHS 890 15 Query: EL K M+ + S

715 EKAELEKSKEMMKKLES 731 Sbict:

Score = 128 (19.2 bits), Expect = 3.9e-04, P = 3.9e-04 Identities = 80/356 (22%), Positives = 148/356 (41%) 20

546 LGESIAANNAYRQQETEHIPRKMPWQSSNHSFPTSIKCLTPHL-----Query: KDGVPGLN-I 597

·LE ++ E+ + +2 + +ZIK L L K 25 G+N +

25 Sbjct: LDEMTQLRDVLETKDKENQTALLEYKSTIHKQEDSIKTLEKELETILSQKKKAEDGINKM 84

598 EELIEKLQSGMVVKDQICD--

VRISDIMDVYEMKLSTLASKESRLQDLLETKALALAQAD 655 30 ++ C M K T + KE + E

85 GKDLFALSREMQAVEENCKNLQKEKDKSNVNHQK-Sbict: ETKSLKEDIAAKITEIKAIN-ENLE 142

35 Query: 656

RLIAGHRCGRTGAETEARTLASMLREVERKNEELSVLLKAGGVESERAGSDIEHLFGHNR 715 ++ Q C EE ++ LE+++ L+

40

Sbict: 743 KWKIG--CNNTZKEKEH--ISKELVEYKSRFQSHDNLVAKLTEKLKSLANNYKDMQAENE 198

716 KLESVAEEHEILTKSYMELLQRN-ESTEKKNKDLQITCDSLNKQIETVKKLNESLKEQNE 774

45 EE + + + LQ +S ++ ++ QI Z+ K IE +KK Sbjct: 199 SLĪKAVEESKNESSIQLSNLQNKIDSMSQEKENFQIERGSIEKNIEQLKKTISDLEQTKE 258

50 Query: 775 KSIAQLIEKEEQRKEVQNQLVDREHKLANLHQKTKVQEEKIKTLQKEREDKEETI---- 829 + I++ + + E ++Q+ + KL KI LK RE+ F

Sbjct: 524 EII2K---

SDSSKDEYESQISLLKEKLETATTANDENVNKISELTKTREELEAELAAYKN 315 55

A30 --DILRKELSRTEQIRKELSIKASSLEVQKAQLEGRLEE-KESLVKLQQ--Querv: EEFNK-HZH 993

WO 01/98454 PCT/IB01/02050 L+ +K QLE + L +L +E+ KE+ E K+ L L+ LKH Sbjct: 37P LKNELETKLETSEKALKEVKENEEHLKEEKIQLEKEATETKQQLNSLRANLESLEKEHED 375 888 IMAIN 488 Query: + A + Sbjct: 376 LAAQL 380 10 Score = 117 (17.6 bits): Expect = 3.8e-03: P = 3.8e-03 Identities = 50/240 (20%), Positives = 111/240 (46%) ASKESRLQDLLETKALALAQADRLIAQHRCQRTQAETEARTLASMLREVERKNEELSVLL 693 15 A E L+ L E + A+ ++ + + E+ L S + + + +E+L Sbjct: 699 AKVEEGLKKLEEESSKEKAELEKSKEMMKKLESTIESNETELKSSMETIRKSDEKLEQSK 758 20 694 KARRVESERAR---SD-IEHLFQHNRKLESVAEEHEILTKSYMELLQRNESTEKKNKDLQ 749 K++++Q SD I + + + +E + + I KZ Et. Sbjct: 759 25 KSAEEDIKNLQHEKSDLISRINESEKDIEELKSKLRIEAKSSSELETVKQELNNAQEKIR AJA ITCDSLNKQIETVKKLNESLKEQNEKSIAQLIEKEEQRKEVQNQLVDREHKLANLHQKTK 809 + + N +++ KL + +E +K A++ +E+++ + ++L + E +| 30 Sbjct: 819 VNAEE-NTVLKS--KLEDIERELKDKQ-AEIKSNGEEKELLTSRLKELEGELDSTGGKAG 874 810 VQEEK----Query: 35 IKTLQKEREDKEETIDILRKELSRTEQIRKELSIKASSLEVQKAQLEGRLE 865 ++ Q E+ +E +L E + + KE + K EE+ + + + Sbjct: 875 KSEEESRAEVRKFQVEKSQLDEKAMLL--ETKYNDLVNKEQAWKRDEDTVKKTT-DSQRQ 931 40 8PP EKEZTAK 945 Query: EELK BEP NALNAIB SEP Sbjct: Score = 109 (16.4 bits), Expect = 2.6e-02, P = 2.5e-02 45 Identities = 64/284 (22%), Positives = 135/284 (47%) Query: 598 EELIEKLQSGMVVKDQICDVRISDIMDVYEMKLSTLASKESRLQDLLETKALALA---QA 654 50 +E+++KL+S + + + I Ε + S E +++L ++ Sbjct: 723 KEMMKKLESTIESNETELKSSMETIRKSDEKLEQSKKSAEEDIKNLQHEKSDLISRINES 782

55 Query: 655 DRLIAQHRCQRTQAETEARTLASMLREVERKNEELSVLLKAQQVESERAQSDIEH-LFQ 712
++ I + + + R +A++ + L ++ +E+ E++ V + V +
DIE L

Sbjct: 783 EKDIEELKSKLRIEAKSSSE-LETVKQELNNAQEKIRVNAEENTVLKSKLE-DIERELKD 840

Query: 713 HNRKLESVAEEHEILTKSYMELLQRNESTEKK-NKDLQITCDSLNK-5 QIETVKKLNES-- 768

+++S EE E+LT EL Q +ST++K K + + + K Q+E

+L+E

Sbjct: 841 KQAEIKSNQEEKELLTSRLKELEGELDSTQQKAQKSEEESRAEVRKFQVEK-SQLDEKAM 899

10

Query: 769 LKEQNEKSIA---QLIEKEEQ--RKEVQNQLVDREHKLANLHQKTKVQEEKIKTLQKERE 823 L E + Q +++E +K +Q + E KLA

K +

+ L K

K+K ++R

15 Sbjct: 900 LLETKYNDLVNKEQAWKRDEDTVKKTTDSQRQEIE-KLAKELDNLKAENSKLKEANEDRS 958

Query: 824 DKEETI----DILRKELSRTEQIRKELSIKASSLEVQKAQLEGRLEEKE

20 + ++ + D+ K ++ K+L ++ SS E + E E+ E
Sbjct: 959 EIDDLMLLVTDLDEKNAKYRSKL-KDLGVEISSDEEDDEEEEDDEEDDE
1006

Score = 96 (14.4 bits), Expect = 1.1e+00, P = 6.6e-01 25 Identities = 40/210 (19%), Positives = 101/210 (48%)

30 +L +
Sbict: 15

Sbjct: 15
ETELKNVRDSLDEMT@LRDVLETKDKEN@TALLEYKSTIHK@EDSIKTLEKELETILS@K 74

Query: 739 ESTE----

35 KKNKDLQITCDSLNKQIETVKKLNESLKEQNEKSIAQLIEKEEQRKEVQNQL 794 + E K KDL +L+++++ V++ ++L+++ +KS + +++ K ++ +

Sbjct: 75 KKAEDGINKMGKDLF----ALSREMQAVEENCKNLQKEKDKSN---VNHQKETKSLKEDI 127

40

Query: 795 VDREHKLANLHQKTKVQEEKIKTLQKERED-KEETIDILRKELSRTEQIRKELSIKASSL 853

+ ++ +++ + + + L KE+E +E ++ + 2 + K

L+ K SL
45 Sbjct: 128 AAKITEIKAINENLEKMKIQCNNLSKEKEHISKELVEYKSRFQSHDNLVAK-LTEKLKSL 186

Query: 854 EVQKAQLEGRLEEKESLVKLQQEELNKHSHMIAMIHS 890

50 Spjct: JB7 ANNYKDMQA---ENESLIKAVEESKNESSIQLSNLQN 220

Score = 52 (7.8 bits), Expect = 2.0e-10, P = 2.0e-10 Identities = 39/167 (23%), Positives = 74/167 (44%)

55 Query: 99 LNSVLAGVVCRSSHTDSVFLQCIQLLQKLTYNVKIFYSGANIDEL-ITFLIDHIQSSEDE 157 LN + +++ ++ L+ I+ ++ T +K N E ++ L D +++SED+

PCT/IB01/02050 WO 01/98454 Sbict: LNLQIKELKKKNETNEASLLESIKSIESETVKIKELQDECNFKEKEVSELEDKLKASEDK 506 158 -Query: LKMPCLGLLANLCRHNLSVQTHIKTLSNVKSFYRTLITLLAHSSLTVVVFALSILSSLT 216 K L + + L +T T ++ T ++ Sbjct: 507 NSKYLEL@KESEKIKEELDAKT---TELKIQLEKVTNLSKAKEKSESELSRLKKTSSEER 563 10 217 LN-EEVGEKLFHARNI-HQTFQLIFNILINGDGTLTRKYS--VDLLMDLL Query: 565 I +Q F+ +L 'G T+T++YS ++ L D L N EE EKL + 564 KNAEEQLEKLKNEIQIKNQAFEKERKLLNEGSSTITQEYSEKINTLEDEL Sbjct: 15 **P13** Pedant information for DKFZphmel2_12jl, frame 2 20 Report for DKFZphmel2 12j1.2 **ELENGTHD** 905 25 EMMI 102067-81 [[q] 5.85 TREMBL:SCINTANA_1 Saccharomyces cerevisiae EHOMOLI integrin analogue gene, complete cds. Le-14 **EFUNCATI** OB.07 vesicular transport (golgi network, etc.) 30 cerevisiae, YDLO58wJ 5e-16 ### CFUNCATI 30-03 organization of cytoplasm ES. cerevisiae, YDL058wl 5e-lb EFUNCATD 1 genome replication, transcription, recombination and repair EM. jannaschii MJ13221 le-10 35 [FUNCAT] 09-10 nuclear biogenesis ES. cerevisiae, YDR356w1 2e-10 [FUNCAT] 3D-04 organization of cytoskeleton ES. cerevisiae, YDR356w3 2e-10 [FUNCAT] 03-22 cell cycle control and mitosis IS cerevisiae. 40 YDR356wl 2e-10 EFUNCATE 30-10 nuclear organization CS- cerevisiae YKRD95wll le-09 EFUNCATI 11.04 dna repair (direct repair, base excision repair and nucleotide excision repair) ES. cerevisiae, YKR095wl le-09 45 EFUNCATD D8-22 cytoskeleton-dependent transport ES cerevisiae; YHRO23w MYOL - myosin-l isoforml 4e-09 EFUNCATE 03.04 budding, cell polarity and filament formation

45 LFUNCATU UB-22 cytoskeleton-dependent transport ES- cerevisiae YHRU23w MYOL - myosin-L isoformU 4e-U9

EFUNCATU UB-04 budding, cell polarity and filament formation

ES- cerevisiae, YHRU23w MYOL - myosin-L isoformU 4e-U9

EFUNCATU UB-25 cytokinesis ES- cerevisiae, YHRU23w MYOL
50 myosin-L isoformU 4e-U9

EFUNCATU 99 unclassified proteins ES- cerevisiae, YNLU91wU

EFUNCATI 09.25 vacuolar and lysosomal biogenesis ES. cerevisiae, YOR326wl 6e-08

55 EFUNCATO DA.Lb extracellular transport ES. cerevisiaen YOR326wD 6e-D8 EFUNCATO D9.13 biogenesis of chromosome structure ES. cerevisiaen YLRO86wD 8e-D8

98 classification not yet clear-cut **EFUNCATI** ES. cerevisiae, YJR134cI le-07 **EFUNCATI** Ob.O7 protein modification (glycolsylation, acylation, myristylation, palmitylation, farnesylation and processing) ES- cerevisiae, YKL201c1 4e-07 30.05 organization of centrosome **EFUNCATI** ES. cerevisiae. YIL144w1 4e-06 EFUNCATD 03.07 pheromone response, mating-type determination, 10 **EFUNCATI** D8.99 other intracellular-transport activities cerevisiae, YNLO79c3 5e-06 09.04 biogenesis of cytoskeleton **EFUNCATI** ES. cerevisiae, YKL179c1 be-06 30.02 organization of plasma membrane 15 **EFUNCATI** ES. cerevisiae; YEROD8c3 8e-Ob **EFUNCATI** 03-19 recombination and dna repair ES. cerevisiae. YNL250wl le-05 03.13 meiosis ES. cerevisiae, YDR285wl le-05 **EFUNCATI** 20 30.13 organization of chromosome structure **EFUNCATI** cerevisiae, YDR285w3 le-05 EFUNCATI 11.01 stress response ES. cerevisiae, YPR141c1 2e-05 **EFUNCATI** Ob.10 assembly of protein complexes II. cerevisiae. YPRl4lc] 2e-05 EFUNCATI 06.01 protein folding and stabilization LZcerevisiae, YNL227cJ 9e-05 EFUNCATI 05-04 translation (initiation, elongation and termination) ES. cerevisiae, YALO35wl le-04
EFUNCATU 10.05.99 other pheromone response activities EZcerevisiae, YHR158cl le-04 30 [FUNCAT] o chaperones EM. genitalium, MG3551 2e-04 [FUNCAT] 03.22.01 cell cycle check point proteins cerevisiae, YGLO86w3 2e-04 EFUNCATI 03.10 sporulation and germination ES. cerevisiae. YNL225cJ 3e-D4 35 **EFUNCATE** r general function prediction EM jannaschii MJ12543 4e-04 [FUNCAT] OB-Ol nuclear transport [S. cerevisiae, YPL174c] 4e-O4 [FUNCAT] 04-05-01-01 general transcription activities 40 EBFOCK21 bk07005E **EBFOCK21** BLO1160B Kinesin light chain repeat proteins EBLOCKSI BLOO326D Tropomyosins proteins d2tmab_ 1.105.4.1.1 Tropomyosin Erabbit EZC0P3 45 (Oryctolagus cuniculus) 3e-23 3.6.1.32 Myosin ATPase 4e-10 **EECJ EPIRKUJ** nucleus 5e-09 **EPIRKWI** phosphotransferase 2e-07 blocked amino end le-Ob **EPIRKU** 50 [PIRKW] duplication 2e-07 citrulline 3e-08 [PIRKW] **IPIRKU**I tandem repeat 4e-10 **EPIRKU**I heterodimer le-0? heart 4e-08 [PIRKW] 55 endocytosis 7e-08 **EPIRKUI** transmembrane protein le-14 **EPIRKWI** serine/threonine-specific protein kinase 2e-07 **EPIRKWI**

cell wall Ze-Ob

EPIRKU

```
WO 01/98454
                                                        PCT/IB01/02050
    EPIRKUI
                    zinc finger 7e-08
    EPIRKU
                    DNA binding 3e-09
    EPIRKU
                    metal binding 7e-08
    CPIRKUJ
                    muscle contraction 4e-10
    EPIRKU
                    brain 2e-06
    EPIRKWI
                    acetylated amino end 2e-07
    EPIRKWI
                    heterotetramer 5e-07
    EPIRKW
                    actin binding 4e-10
    EPIRKWI
                    mitosis le-O8
10
    EPIRKWI
                    microtubule binding le-D8
    EPIRKU
                  · ATP 4e-10
    EPIRKWI
                    chromosomal protein le-D?
                    thick filament 9e-10
    TPIRKW
    IPIRKWI
                    phosphoprotein le-09
15
    EPIRKWI
                    skeletal muscle le-O8
                    calcium binding 3e-08
    EPIRKWI
    [PIRKW]
                    alternative splicing 9e-10
                    DNA condensation le-07
    CPIRKWI
    CPIRKW]
                    coiled coil Le-14
20
    EPIRKU
                    P-loop 2e-10
    EPIRKWI
                   heptad repeat 5e-09
    EPIRKUJ
                   methylated amino acid 4e-10
    EPIRKWI
                    immunoglobulin receptor 2e-07
    EPIRKWI
                   peripheral membrane protein 7e-08
25
    EPIRKWI
                   cardiac muscle 4e-DA
    EPIRKU
                   hydrolase 4e-10
    EPIRKWI
                   microtubule 5e-09
    EPIRKWI
                   muscle 4e-D8
    EPIRKUJ
                   membrane protein 5e-09
30
    EPIRKUJ
                   EF hand 3e-08
    EPIRKWI
                   cell division le-Ob
    EPIRKUJ
                   cytoskeleton be-09
    EPIRKUJ
                   hair 3e-D8
    EPIRKWI
                   calmodulin binding 7e-08
    EPIRKUJ
                   Golgi apparatus 2e-07
    ESUPFAMI
              hypothetical protein YJL074c 5e-09
    ESUPFAM3
              unassigned Ser/Thr or Tyr-specific protein kinases 2e-
    07
    ESUPFAMI
              myosin motor domain homology 2e-10
40
    ESUPFAM3
              alpha-actinin actin-binding domain homology be-D9
    ESUPFAMI
              tropomyosin 2e-08
    ESUPFAM3
              kinesin heavy chain 5e-07
    ESUPFAM3
              plectin be-09
              SAM homology le-Ob
    ESUPFAM3
    ESUPFAM3
              trichohyalin 3e-08
    ESUPFAM3
              ribosomal protein S10 homology be-09
    ESUPFAM3
              protein kinase C zinc-binding repeat homology 5e-09
    ESUPFAME
              giantin 7e-08
    ESUPFAMI
              protein kinase homology 2e-07
   ESUPFAMI
              protein 4-1 membrane-binding domain homology 9e-08
    ESUPFAMI
             human early endosome antigen 1 7e-08 myosin MY02 2e-06
    ESUPFAMI
   ESUPFAMI
             M5 protein 3e-09
   ESUPFAM3
             Mycoplasma genitalium hypothetical protein MG218 5e-09
   ESUPFAMD
             myosin heavy chain 2e-10
             conserved hypothetical P115 protein 3e-09
   ESUPFAMD
             centromere protein E le-D&
   ESUPFAMI
             calmodulin repeat homology 3e-08
   ESUPFAMI
```

WO 01/98454 PCT/IB01/02050 **ESUPFAMJ** hypothetical protein MJ0914 2e-07 ESUPFAMD hypothetical protein MJ1322 3e-09 **ESUPFAMD** pleckstrin repeat homology 5e-D9 kinesin motor domain homology le-D& ESUPFAMI 5 [CMA79U2] ezrin "Te=88 **IPROSITED LEUCINE_ZIPPER L** [KW] TRANSMEMBRANE [KW] LOW_COMPLEXITY 3.09 % EKUI COILED_COIL 18.34 % 10 MDSTACLKSLLLTVSQYKAVKSEANATQLLRHLEVISGQKLTRLFTSNQILTSECLSCLV SEQ SEG PRD COILS 15 MEM SEQ ELLEDPNISASLILSIIGLLSQLAVDIETRDCLQNTYNLNSVLAGVVCRSSHTDSVFLQC 20 SEG PRD COILZ MEM ----- MMMMMMMMMMMMMM ------25 IQLLQKLTYNVKIFYSGANIDELITFLIDHIQSSEDELKMPCLGLLANLCRHNLSVQTHI ZEQ SEG PRD COILS 30 MEM KTLSNVKSFYRTLITLLAHSSLTVVVFALSILSSLTLNEEVGEKLFHARNIHQTFQLIFN SEQ SEG 35 PRD COILS MEM ILINGDGTLTRKYSVDLLMDLLKNPKIADYLTRYEHFSSCLHQVLGLLNGKDPDSSSKVL 40 SEQ SEG PRD COILS 45 MEM ELLLAFCSVTQLRHMLTQMMFEQSPPGSATLGSHTKCLEPTVALLRWLSQPLDGSENCSV SEQ SEG PRD 50 COILZ MEM SEQ LALELFKEIFEDVIDAANCSSADRFVTLLLPTILDQLQFTEQNLDEALTRKNVKGIAKAI 55 SEG PRD COILS

	W	O 01/98454	PCT/IB01/02050
	MEM	•••••	
	SEQ SEG	EVLLTLCGDDTLKMHIAKILTTVKCTTLIEQQFTY0	SKIDLGFGTKVADSELCKLAADVIL
5	PRD COIL	hhhhhhccccchhhhhhhhhhhhheeeeeeeeecc.S	cccccceeehhhhhhhhhhhhh
	MEM	••••••••••••••	
10	SEQ	KTLDLINKLKPLVPGMEVSFYKILØDPRLITPLAFA	NLTSDNREQVQSGLRILLEAAPLPI
	SEG PRD COIL	hhhhhhhhcccccccccceeeccccchhhhhr	nhccccchhhhhhhhhhhhhcccc
15	MEM	••••••••••••••	•••••••••
	SEQ SEG	FPALVLGESIAANNAYRQQETEHIPRKMPWQSSNHS	:FPTSIKCLTPHLKDGVPGLNIEEL
20	PRD COIL:	cceeeehhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh	ccchhhhhhhhhhhhhhhhhhh
	MEM		••••••••••••••••
25	SEQ	IEKL@SGMVVKD@ICDVRISDIMDVYEMKLSTLASK	
	PRD COIL:	hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh	
30	MEM	••••••••••••	
	SEG SEG	HRCQRTQAETEARTLASMLREVERKNEELSVLLKAQ	
	PRD COILS		
35	MEM	CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC	
	SEQ SEQ	AEEHEILTKSYMELLQRNESTEKKNKDLQITCDSLN	KGIETVKKLNESLKEGNEKSIAGL
40	PRD COILS	հերորդում և անակաների հերերերի հերերերի հերերեր Հ	hhhhhhhhhhhhhhhhhhhhh
	MEM		
45	SEQ SEG	IEKEE@RKEV@N@LVDREHKLANLH@KTKV@EEKIK	TLQKEREDKEETIDILRKELSRTE
		_ հեռերերերեր և Հայաստում և հեռերերեր և հեռերեր և հեռերեր և հեռերերեր և հեռերերեր և հեռերերեր և հեռերեր և հեռեր Տ	հիհիհիհիհիհիհիհիհիհի
50	MEM	······································	
	SEQ	@IRKELSIKASSLEV@KA@LEGRLEEKESLVKL@@EI	ELNKHSHMIAMIHSLSGGKINPET
55	SEG PRD COILS	իհիրինինինինինինինինինինինինինինինի	nhhhhhhhhhhhhhhhhhcccc
		CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC	CCCC
	MEM		

SEQ VNLSI
SEG
PRD CCCCC
COILS
5 MEM

Prosite for DKFZphmel2_12j1.2

10 PS00029

331->353 LEUCINE_ZIPPER

PS000029

(No Pfam data available for DKFZphmel2_12j1.2)
15

5 group: intracellular transport and trafficing

DKFZphmel2_7g14 encodes a novel 973 amino acid protein with similarity to the dor (deep orange) protein of drosophila melanogaster.

The novel protein is also similar to the vakuolar membrane protein pep3 of Saccharomyces cerevisiae, which is involved in protein sorting mechanisms. The expression profile is ubiquitous and a role in protein transport/targeting is likely.

The new protein can find application in modulation of the sorting of proteins into different compartments.

20 similarity to DEEP ORANGE (Drosophila melanogaster)

perhaps complete cds- and full length

Sequenced by MediGenomix

25 Locus: unknown

30

Insert length: 3951 bp

Poly A stretch at pos. 3893, polyadenylation signal at pos. 3874

1 GCCCGCGTCA CGGGGGCGGG AGTCAGCTGA GCTGCCGGGG CGAGGTTGGG 51 ATCACCTGGC ACCGGCTGAA GGGAGCCTGT GATTTTTTTG TAGCGGGGGC LOL GGGGAGTAAG GTGCAAGACT GCGCCAGATT CAAGGACGAG GGCTGCCCGA 151 TTATCTCGCT GCATAAGGCA. AGAGCAAGAG. GATCCTCAGG ATTTTAAAGA 35 201 GGAGGCGACG GCTGCAGGTT CCCAGGATCT GTCAGAGGCT GGGGAGTTAC 251 AGCTTCCATT CTGGGGCGAC GGGGGACCCG GGGGGGTAGC CCTTTTGTAA
301 TCCCCAGGCC CCGGACAAAG AGCCCAGAGG CCGGGCACCA TGGCGTCCAT
351 CCTGGATGAG TACGAGAACT CGCTGTCCCG CTCGGCCGTC TTGCAGCCCG
401 GCTGCCCTAG CGTGGGCATC CCCCACTCGG GGTATGTGAA TGCCCAGCTG
451 GAGAAGGAAG TGCCCATCTT CACAAAGCAG CGCATTGACT TCACCCCTTC
501 CGAGCGCATT ACCAGTCTTG TCGTCTCCAG CAATCAGCTG TGCATGAGCC
551 TGGGCAAGGA TACACTGCTC CGCATTGACT TGGGCAAGGC AAATGAGCCC 40 LOI AACCACGTGG AGCTGGGACG TAAGGATGAC GCAAAAGTTC ACAAGATGTT
LSI CCTTGACCAT ACTGGCTCTC ACCTGCTGAT TGCCCTGAGC AGCACGGAGG
701 TCCTCTACGT GAACCGAAAT GGACAGAAGG TACGGCCACT AGCACGCTGG 45 751 AAGGGCAGC TGGTGGAGAG TGTGGGTTGG AACAAGGCAC TGGGCACGGA BD1 GAGCAGCACA GGCCCCATCC TGGTCGGGAC TGCCCAAGGC CACATCTTTG B51 AAGCAGAGCT CTCAGCCAGC GAAGGTGGGC TTTTCGGCCC TGCTCCGGAT 50 901 CTCTACTTCC GCCCATTGTA CGTGCTAAAT GAAGAAGGGG GTCCAGCACC 951 TGTGTGCTCC CTTGAGGCCG AGCGGGGCCC TGATGGGCGT AGCTTTGTTA BODD TTGCCACCAC TCGGCAGCGC CTCTTCCAGT TCATAGGCCG AGCAGCAGAG 1051 GGGGCTGAGG CCCAGGGTTT CTCAGGGCTC TTTGCAGCTT ACACGGACCA 1101 CCCACCCCA TTCCGTGAGT TTCCCAGCAA CCTGGGCTAC AGTGAGTTGG 1151 CCTTCTACAC CCCCAAGCTG CGCTCCGCAC CCCGGGCCTT CGCCTGGATG 55 1201 ATGGGGGATG GTGTGTTGTA TGGGGCATTG GACTGTGGGC GCCCTGACTC 1251 TCTGCTGAGC GAGGAGCGAG TCTGGGAGTA CCCAGAGGGG GTAGGGCCTG ₽30P GGGCCAGCCC ACCCCTAGCC ATCGTCTTGA CCCAGTTCCA CTTCCTGCTG

```
1351 CTACTGGCAG ACCGGGTGGA GGCAGTGTGC ACACTGACCG GGCAGGTGGT
        1401 GCTGCGGGAT CACTTCCTGG AGAAATTTGG GCCGCTGAAG CACATGGTGA
        1451 AGGACTCCTC CACAGGCCAG CTGTGGGCCT ACACTGAGCG GGCTGTCTTC
1501 CGCTACCACG TGCAACGGGA GGCCCGAGAT GTCTGGCGCA CCTATCTGGA
1551 CATGAACCGC TTCGATCTGG CCAAAGAGTA TTGTCGAGAG CGGCCCGACT
  5
        1601 GCCTGGACAC GGTCCTGGCC CGGGAGGCCG ATTTCTGCTT TCGCCAGCGT
        1651 CGCTACCTGG AGAGCGCACG CTGCTATGCC CTGACCCAGA GCTACTTTGA
        1701 GGAGATTGCC CTCAAGTTCC TGGAGGCCCG ACAGGAGGAG GCTCTGGCTG
        1751 AGTTCCTGCA GCGAAAACTG GCCAGTTTGA AGCCAGCCGA ACGTACCCAG
        1801 GCCACACTGC TGACCACCTG GCTGACAGAG CTCTACCTGA GCCGGCTTGG
10
        1851 GGCTCTGCAG GGCGACCCAG AGGCCCTGAC TCTCTACCGA GAAACCAAGG
        1901 AATGCTTTCG AACCTTCCTC AGCAGCCCCC GCCACAAGA GTGGCTCTTT
        1951 GCCAGCCGGG CCTCTATCCA TGAGCTGCTC GCCAGTCATG GGGACACAGA
        2001 ACACATGGTG TACTTTGCAG TGATCATGCA GGACTATGAG CGGGTGGTGG
       2051 CTTACCACTG TCAGCACGAG GCCTACGAGG AGGCCCTGGC CGTGCTCGCC
15
    2051 CTTACCACTG TCAGCACGAG GCCTACGAGG AGGCCCTGGC CGTGCTCGCC
2101 CGCCACCGTG ACCCCCAGCT CTTCTACAAG TTCTCACCCA TCCTCATCCG
2151 TCACATCCCC CGCCAGCTTG TAGATGCCTG GATTGAGATG GGCAGCCGGC
2201 TGGATGCTCG TCAGCTCATT CCTGCCCTGG TGAACTACAG CCAGGGTGGT
2251 GAGGTCCAGC AGGTGAGCCA GGCCATCCGC TACATGGAGT TCTGCGTGAA
2301 CGTGCTGGGG GAGACTGAGC AGGCCATCCA CAACTACCTG CTGTCACTGT
2351 ATGCCCGTGG CCGGCCGGAC TCACTACTGG CCTATCTGGA GCAGGCTGGG
2401 GCCAGCCCCC ACCGGGTGCA TTACGACCTC AAGTATGCGC TGCGGCTCTG
2451 CGCCGAGCAT GGCCACCACC GCGCTTGTGT CCATGTCTAC AAGGTCCTAG
2501 AGCTGTATGA GGAGGCCGTG GACCTGGCC TGCAGGTGGA TGTGGACCTG
2551 GCCAAGCAGT GTGCAGACCT GCCTGAGGAG GATGAGGAAT TGCGCAAGAA
2501 GCTGTGGCTG AAGATCGCAC GGCACGTGGT GCAGGAAGAG GAAGATGTAC
2551 AGACAGCCAT GGCTTGCCTG GCTAGCTGC CCTTGCTCAA GATTGAGGAT
20
25
       2651 AGACAGCCAT GGCTTGCCTG GCTAGCTGCC CCTTGCTCAA GATTGAGGAT
       2701 GTGCTGCCCT TCTTTCCTGA TTTCGTCACC ATCGACCACT TCAAGGAGGC
       2751 GATCTGCAGC TCACTTAAGG CCTACAACCA CCACATCCAG GAGCTGCAGC
       2801 GGGAGATGGA AGAGGCTACA GCCAGTGCCC AGCGCATCCG GCGAGACCTG
30
       2851 CAGGAGCTGC GGGGCCGCTA CGGCACTGTG GAGCCCCAGG ACAAATGTGC
       2901 CACCTGCGAC TTCCCCCTGC TCAACCGCCC TTTTTACCTC TTCCTCTGTG
       2951 GCCATATGTT CCATGCTGAC TGCCTGCTGC AGGCTGTGCG ACCTGGCCTG
       BDDL CCAGCCTACA AGCAGGCCCG GCTGGAGGAG CTGCAGAGGA AGCTGGGGGC
       3051 TGCTCCACCC CCAGCCAAGG GCTCTGCCCG GGCCAAGGAG GCCGAGGGTG
35
       3101 GGGCTGCCAC GGCAGGGCCC AGCCGGGAAC AGCTCAAGGC TGACCTGGAT 3151 GAGTTGGTGG CCGCTGAGTG TGTGTACTGT GGGGAGCTGA TGATCCGCTC 3201 TATCGACCGG CCGTTCATCG ACCCCCAGCG CTACGAGGAG GAGCAGCTCA
       BESL GTTGGCTGTA GGAGGGTGTC ACCTTTGATG GGGGTGGGCA ATGGGGAGCA
40
       BELL BAGGCTTGAA CCCACTTGAG AAGGCTGCCT CCTAGGCTCT GCTCAGTCAT
       3351 CTTGCAATTG CCACACTGTG ACCACGTTGA CGGGAGTAGA GTAGCGCTGT
       3401 TGGCCAGGAG GTGTCAGGTG TGAGTGTATT CTGCCAGCTT TTCATGCTGT
       3451 TCTTCAGAGC TGCAGTTATG CCAGACCATC AGCCTGCCTC CCAGTAGAGG
       1501 CCCTTCACCT GGAGAAGTCA GAAATCTGAC CCCATTCCAC CCCCTGCCTC
    3551 TAGCACCTCT TCTGTCCCTG TCATTCCCCA CACACGTCCT GTTCACCTCG
45
       3601 AGAGAGAGA AGAGAGAGCA CCTTTCTTCC GTCTGTTCAC TCTGCGGCCT
       3651 CTGGAATCCC AGCTCTTCTC TCTCAGAAGA AGCCTTCTCT TCCTCCTGCC
       3701 TGTAGGTGTC CCAGAAGTGA GAAGGCAGCC TTCGAAGTCC TGGGCATTGG
       3751 GTGAGAAAGT GATGCTAGTT GGGGCATGCT TTTGTGCACA CTCTCTGGGG
       3801 CTCCAGTGTG AAGGGTGCCC TGGGGCTGAG GGCCTTGTGG AGGATGGTCG
50
       3851 GTGGTGGTGA TGGAGGTGGA GAGCATTAAA CTGTCTGCAC TGCAAAAAA
       3951 A
```

55

PCT/IB01/02050

No BLAST result

5

Medline entries

97218037:

Shestopal SA, Makunin IV, Belyaeva ES, Ashburner M, Zhimulev IF.;Mol

10 Gen Genet 1997 Feb 20:253(5):642-8

92049306:

Robinson JS, Graham TR, Emr SD.; A putative zinc finger protein, Saccharomyces cerevisiae

15 Vpslap, affects late Golgi functions required for vacuolar protein sorting and efficient alpha-factor prohormone maturation. Mol Cell Biol 1991 Dec;11(12):5813-24

92049305:

20 Preston RA, Manolson MF, Becherer K, Weidenhammer E, Kirkpatrick

Wright Ra

Jones EW-; Isolation and characterization of PEP3, a gene required

25 for vacuolar biogenesis in Saccharomyces cerevisiae. Mol Cell Biol 1991
Dec:11(12):5801-12

30

Peptide information for frame 1

35 ORF from 340 bp to 3258 bp; peptide length: 973 Category: similarity to known protein Classification: Cellular transport and traffic

1 MASILDEYEN SLSRSAVL@P GCPSVGIPHS GYVNA@LEKE VPIFTK@RID 40 51 FTPSERITSL VVSSNQLCMS LGKDTLLRID LGKANEPNHV ELGRKDDAKV IDI HKMFLDHTGS HLLIALSSTE VLYVNRNGQK VRPLARWKGQ LVESVGWNKA 151 LGTESSTGPI LVGTAQGHIF EAELSASEGG LFGPAPDLYF RPLYVLNEEG 201 GPAPVCSLEA ERGPDGRSFV IATTRQRLFQ FIGRAAEGAE AQGFSGLFAA 251 YTDHPPPFRE FPSNLGYSEL AFYTPKLRSA PRAFAWMMGD GVLYGALDCG 301 RPDSLLSEER VWEYPEGVGP GASPPLAIVL TQFHFLLLLA DRVEAVCTLT 45 351 GQVVLRDHFL EKFGPLKHMV KDSSTGQLWA YTERAVFRYH VQREARDVWR 403 TYLDMNRFDL AKEYCRERPD CLDTVLAREA DFCFRQRRYL ESARCYALTQ 451 SYFEEIALKF LEARGEEALA EFLQRKLASL KPAERTQATL LTTULTELYL 501 SRLGALQGDP EALTLYRETK ECFRTFLSSP RHKEWLFASR ASIHELLASH 551 GDTEHMVYFA VIMQDYERVV AYHCQHEAYE EALAVLARHR DPQLFYKFSP 50 LOT ILIRHIPROL VDAWIEMGSR LDAROLIPAL VNYSOGGEVO OVSQAIRYME LITRHIPROL VDAWIENGSK LDAROLIPAL VNYSQGGEVO WYSQAIRYNE LSI FCVNVLGETE QAIHNYLLSL YARGRPDSLL AYLEQAGASP HRVHYDLKYA 701 LRLCAEHGHH RACVHVYKVL ELYEEAVDLA LQVDVDLAKQ CADLPEEDEE 751 LRKKLWLKIA RHVVQEEEDV QTAMACLASC PLLKIEDVLP FFPDFVTIDH 801 FKEAICSSLK AYNHHIQELQ REMEEATASA QRIRRDLQEL RGRYGTVEPQ 851 DKCATCDFPL LNRPFYLFLC GHMFHADCLL QAVRPGLPAY KQARLEELQR 901 KLGAAPPPAK GSARAKEAEG GAATAGPSRE QLKADLDELV AAECVYCGEL 951 MIRSIDRPFI DPQRYEEEQL SWL 55

BLASTP hits

5 No BLASTP hits available

Alert BLASTP hits for DKFZphmel2_7gl4, frame l

- 10 SWISSPROT:DOR_DROME DEEP ORANGE PROTEIN., N = 1, Score = 1279, P = 2.4e-130
- PIR:A41943 vacuolar membrane protein PEP3 yeast (Saccharomyces cerevisiae), N = 3, Score = 266, P = 5.1e-27

>SWISSPROT:DOR_DROME DEEP ORANGE PROTEIN.

Length = 1.002

20

HSPs:

Score = 1279 (191.9 bits), Expect = 2.4e-130, P = 2.4e-130 Identities = 303/847 (35%), Positives = 463/847 (54%)

25
Query: 130
KVRPLARWKGQLVESVGWNKALGTESSTGPILVGTAQGHIFEAELSASEGGLFGPAPDLY 189
KVR + ++K + +V +N G ESSTGPIL+GT++G IFE EL+ + G

30 Sbjct: 155 KVRRIEKFKDHEITAVAFNPYHGNESSTGPILLGTSRGLIFETELNPAADG-----HVQ 208

Query: 190 FRPLYVLNEEGGPA-PVCSLEAERGPDG-RSFVIATTRQRLFQFIGRAAEGAEAQGFSGL 247

- 35 + LY L G P P+ L+ R P+ R ++ T+ +++ F
 AE + +
 Sbjct: 209 RKQLYDLGL-GRPKYPITGLKLLRVPNSSRYIIVVTSPECIYTF-QETLKAEERSLQAI 265
- 40 Query: 248 FAAYTD-HPPPFREFPSNLGYSELAFYTPKLRSAPAFAWMGDGVLYGAL--DCGRPD 3D3
 FA Y P E ++L +S+L F+ P P+ +W+ G+G+ G L
 +

Sbjct: 266

45 FAGYVSGVQEPHCEERKTDLTFSQLRFFAPPNSKYPKQWAWLCGEGIRVGELSIEANSAA 325

Query: 304 SLLSEERV---WEYPEGVGPGA---SPPLAIVLTQFHFLLLLADRVEAVCTLTGQVVLRD 357

+L+ + +E + G + P A VLT++H +LL AD V A+C L

- 50 + V ++
 Sbjct: 326
 TLIGNTLINLDFEKTMHLSYGERRLNTPKAFVLTEYHAVLLYADHVRAICLLNGEQVYQE 385
- Query: 358 HFLE55 KFGPLKHMVKDSSTGQLWAYTERAVFRYHVQREARDVWRTYLDMNRFDLAKEYCR 416
 FE+G++DTG++YT+VFVRER+VWRYLD
 +++LA+

Sbjct: 386

AFDEARVGKPLSIERDELTGSIYVYTVKTVFNLRVTREERNVWRIYLDKGQYELATAHAA 445

Query: 417

ERPDCLDTVLAREADFCFRQRRYLESARCYALTQSYFEEIALKFLEARQEEALAEFLQRK 476 E P+ L VL + AD F Y +A YAT FEE+ LKF+

Sbjct: 446

EDPEHLQLVLCQRADAAFADGSYQVAADYYAETDKSFEEVCLKFMVLPDKRPIINYVKKR 505

10

15

477 LASL--KPAERXXXXXXXXXXXXXXXXRLGALQ----GDPEALTLYRETKEC-FRTFLSS 529 L+ + KP E

Sbict: 506

LSRVTTKPMETDELDEDKMNIIKALVIWLIDLYLIQINMPDKDEEWRSSWQTEYDEFMME 565

LL

P+

+R + +

Query: 530

PRHKEWLFASRASIHELLASHGDTEHMVYFAVIMQDYERVVAYHCQHEAYEEALAVLARH 589 +R ++ +L+A H D +M FA+ + DY+ VVA

20 EAL L

> Sbjct: 566

AHVLSCTRUNRETVRULIAEHADPRNMAQFAIAIGDYDEVVAQQLKAECYAEALQTLINQ 625

25 Query: 590 RDPQLFYKFSPILIRHIPRQLVDAWIEMGSRLDARQLIPALVNYSQGGEVQQVSQAIRYM 649 R+P+LFYK++P LI +P+ VDA + GSRL+ +L+P L+

Sbjct: 626 RNPELFYKYAPELITRLPKPTVDALMAQGSRLEVEKLVPTLI-

30 IMENREQREQTQ--RYL 682

> Query: 650

EFCVNVLGETEQAIHNYLLSLYARGRPDSLLAYLEQAGASPHRVHYDLKYALRLCAEHGH 709 EF + L T AIHN+LL LYA P L+ YLE G VHYD+ YA

35 ++(+

> Sbjct: EBa

EFAIYKLNTTNDAIHNFLLHLYAEHEPKLLMKYLEIQGRDESLVHYDIYYAHKVCTDLDV 742

HRACVHVYKVLELYEEAVDLALQVDVDLAKQCADLPEEDEELRKKLWLKIARHVVQEEED 769 40 A V + +L + AVDLAL D+ LAK+ A P D ++R+KLWL+IA

Sbjct: 743 KEARVFLECMLRKWISAVDLALTFDMKLAKETASRPS-DSKIRRKLWLRIAYHDIKGTND BOL

45

770 Query: VQTAMACLASCPLLKIEDVLPFFPDFVTIDHFKEAICSSLKAYNHHIQELQREMEEATAS 829 V+ A+ L C LL+IED+LPFF DF ID+FKEAIC +L+ YN

IQELQREM E T

50 Sbict: 802 VKKALNLLKECDLLRIEDLLPFFADFEKIDNFKEAICDALRDYNQRIQELQREMAETTEQ 863

Query: 830

AQRIRRDLQELRGRYGTVEPQDKCATCDFPLLNRPFYLFLCGHMFHADCLLQAVRPGLPA 889 55 R +LQ+LR TVE QD C C+ LL +PF++F+CGH FH+DCL +

V P L

Sbjct: VP5

TDRATAELQQLRQHSLTVESQDTCEICEMMLLVKPFFIFICGHKFHSDCLEKHVVPLLTK 921

Query: 890 YKQARLEELQRKLGAAPPPXXXXXXXXXXXXXXXXXXPSREQLKADLDELVAAECVYCGE 949 + RL L+++L A R LK +++++AA+C++CG Sbict: 922 EQCRRLGTLKQQLEAEVQTQAQPQSGALSKQQAMELQRKRAALKTEIEDILAADCLFCG- 980 950 LMIRSIDRPFIDPQRYEEEQLSW 972 Query: L+I +ID+PF+D +E+ + W10 981 LLISTID@PFVDD--WE@VNVEW 1001 Sbjct: Score = 268 (40.2 bits), Expect = 3.6e-19, P = 3.6e-19 Identities = 91/281 (32%), Positives = 146/281 (51%) 15 36 QLEKEVPIFTKQRIDF-TPSE---RITSLVVSSNQLCMSLG---Query: KDTLLRIDLGKANEPN 88 + ++E IF++ ++ PS + L VS N L LG + TLLR L +A P Sbjct: 20 37 ETDEEDEIFSRHKMVLRVPSNCTGDLMHLAVSRNWLVCLLGTPERTTLLRFFLPRAIPPG 96 B9 HVELGRK---DDAKVHKMFLDHTGSHLLIAL---SST-----EVLYVN--Query: RNGQ----KV 131 25 L + K+ +MFLD TG H++IAL T+2 + LY++ ΚV Sbict: EAVLEKYLSGSGYKITRMFLDPTGHHIIIALVPKSATAGVSPDFLYIHCLESPQAQQLKV 156 30 Query: RPLARWKGQLVESVGWNKALGTESSTGPILVGTAQGHIFEAELSASEGGLFGPAPDLYFR 191 R + ++K + +V +N G ESSTGPIL+GT++G IFE EL+ + G Sbjct: 157 RRIEKFKDHEITAVAFNPYHGNESSTGPILLGTSRGLIFETELNPAADG---35 ---HVQRK 230 192 PLYVLNEEGGPA-PVCSLEAERGPDG-RSFVIATTRQRLFQFIGRAAEGAEAQGFSGLFA 249 G P P+ L+ R P+ R ++ T+ +++ F LY L ΑE 40 231 QLYDLGL-GRPKYPITGLKLLRVPNSSRYIIVVTSPECIYTF--Sbict: **QETLKAEERSLQAIFA 267** Query: 250 AYTD--HPPPFREFPSNLGYSELAFYTPKLRSAPRAFAWMMGDGVLYGAL 45 297 E ++L +S+L F+ P P+ +AU+ G+G+ G L 268 GYVZGVQEPHCEERKTDLTFZQLRFFAPPNZKYPKQWAWLCGEGIRVGEL Sbjct: 317 50 Pedant information for DKFZphmel2_7gl4, frame 1 Report for DKFZphme12_7gl4-1 55

ELENGTHI 973 EMWI 110186-09

	Epil 5.72	
	EHOMOLI SWISSPROT: DOR_DROME DEEP ORANGE PROTEIN. Le-145 EFUNCATI 30.25 vacuolar and lysosomal organization ES.	
	LFUNCATI 30-25 vacuolar and lysosomal organization ES. cerevisiae YLR148wl 5e-41	
5	EFUNCATE Ob-04 protein targeting, sorting and translocation	
	ES. cerevisiae, YLR148wl 5e-41	
	[FUNCAT] OA-O7 vesicular transport (golgi network, etc.) [[S.	-
	cerevisiae YLR148wl 5e-41	
10	EBLOCKSI BLOOlObF Galactokinase proteins	
10	EBLOCKZI PROJO94B	
	IBLOCKZI PRODPOB	
, .	<pre>CPIRKUl yeast vacuole le-39</pre>	
	<pre>EPIRKWl transmembrane protein le-39</pre>	
15	EKWI Alpha_Beta	
	<pre>EKW□ LOW_COMPLEXITY 3.39 % EKW□ COILED_COIL 4.83 %</pre>	
	rvma Colffa 4.02 %	
		•
20	SEG MASILDEYENSLSRSAVLQPGCPSVGIPHSGYVNAQLEKEVPIFTKQRIDFTPSERI	
	SEG	
	PRD ccceeeccccceeeeeccccchhhhhhhhhhhhhhhhh	566

25		
	SEG VVSSNQLCMSLGKDTLLRIDLGKANEPNHVELGRKDDAKVHKMFLDHTGSHLLIALS	STE
	PRD eeccceeeecccccceeeechhhhhhhhheeecccccceeeeee	• • •
	COILS	
30		
	DEA DE MINERAL AND A DESCRIPTION OF THE PROPERTY OF THE PROPER	
	SEQ VLYVNRNGQKVRPLARWKGQLVESVGWNKALGTESSTGPILVGTAQGHIFEAELSASE	
	PRD eeeeeccccchhhhhcccceeeeeecccccccccceeeeecccchhhhhh	
35	COILZ	
	••••••••••••	
	SEQ LFGPAPDLYFRPLYVLNEEGGPAPVCSLEAERGPDGRSFVIATTRQRLFQFIGRAAE	- A -
	SEG	JAL
40	PRD cccccccceeeeccccccceeeeccccccceeeeehhhhhh	
	COILZ	
		· · ·
	SEQ AGGFSGLFAAYTDHPPPFREFPSNLGYSELAFYTPKLRSAPRAFAWMMGDGVLYGALI	N C C
45	SEG	
	PRD hhhchhhhhhhcccccccccccccceeeeccccchhhhhh	ccc
	COILZ	
	••••••	• • •
50	SEQ RPDSLLSEERVWEYPEGVGPGASPPLAIVLTQFHFLLLLADRVEAVCTLTGQVVLRD	JE1
	SEG	
	PRD ccccchhhhhhhcccccccccchhhhhhhhhhhhhhh	ıhh
	COILZ	
55	•••••••••••••••••••••••••••••••••••••••	• • •
	SEQ EKFGPLKHMVKDSSTGQLWAYTERAVFRYHVQREARDVWRTYLDMNRFDLAKEYCREF	ים פ
	SEG	
	PRD hcccccccccccceeeehhhhhhhhhhhhhhhhhhhhh	

COILS CLDTVLAREADFCFRQRRYLESARCYALTQSYFEEIALKFLEARQEEALAEFLQRKLASL SEQ 5 SEG PRD COILZ 10 SEQ KPAERTQATLLTTWLTELYLSRLGALQGDPEALTLYRETKECFRTFLSSPRHKEWLFASR. SEG PRD COILS 15 **ASIHELLASHGDTEHMVYFAVIMQDYERVVAYHCQHEAYEEALAVLARHRDPQLFYKFSP** SEQ PRD COILZ 20 SEQ ILIRHIPRQLVDAWIEMGSRLDARQLIPALVNYSQGGEVQQVSQAIRYMEFCVNVLGETE SEG PRD 25 COILZ QAIHNYLLSLYARGRPDSLLAYLEQAGASPHRVHYDLKYALRLCAEHGHHRACVHVYKVL . SEQ SEG 30 hhhhhhhhhhhhcccccccchhhhhhhhhhhhcccccceeehhhh PRD COILS SEQ ELYEEAVDLALQVDVDLAKQCADLPEEDEELRKKLWLKIARHVVQEEEDVQTAMACLASC 35 SEG PRD COILS PLLKIEDVLPFFPDFVTIDHFKEAICSSLKAYNHHIQELQREMEEATASAQRIRRDLQEL 40 SEQ SEG PRD COILZ 45 RGRYGTVEPQDKCATCDFPLLNRPFYLFLCGHMFHADCLLQAVRPGLPAYKQARLEELQR SEQ SEG PRD COILZ 50 KLGAAPPPAKGSARAKEAEGGAATAGPSREQLKADLDELVAAECVYCGELMIRSIDRPFI SEQ SEG PRD 55 COILS

PCT/IB01/02050

WO 01/98454

SEQ DPQRYEEEQLSWL

PRD chhhhhhhhhccc

5

(No Prosite data available for DKFZphmel2_7gl4.l)

(No Pfam data available for DKFZphmel2_7gl4.l)

5 group: melanoma derived

DKFZphmel2_7kl9 encodes a novel 234 amino acid protein without similarity to known proteins.

- Transcpripts can be found in almost any tissue, but are most abundant in kidney and retina.

 No informative BLAST results: No predictive prosite, pfam or SCOP motife.
- 15 The new protein can find application in studying the expression profile of melanoma-specific genes.

unknown protein

first ATG in frame 1

Sequenced by MediGenomix

25 Locus: /map="3"

20

Insert length: 2386 bp
Poly A stretch at pos. 2343, polyadenylation signal at pos. 2323



1401 GGTAAAATTA AAGAGGAAAT TAATCTTTAT ATATTATTTC TTGCAGAAAC 1451 ATTCATTATT TTATTAATAT TGCCCTAAGT ACAACTAGGC AAGTGATTGC 1501 CACCTAAATC AGAAGACGTT CTAAAGTCAG TAAGAAAGTG TGAAATGCTA 1551 GTATAAAGGT TATTTTTTT CTTTCCTAAA TAACTAAAGT GAGGTGTAGA 160 TTGAGCCTTG ATATTATTTA GTTAATGTTT TTTATTAATT AATTTTGGCT 5 165 GGACTTTATT TAGCTTGATT AGGTTATTAT CTGTCAAACC TTTTAAGTTG 1701 ACAACATGAC TCATATATA ACATGTGTAT AAGATGAGCA TGTGTCGAAG 1751 ACTTATTCGA CTCATTAATG AGGAAACCAG CAGATAGTAA ACCTGGTTCA LBOL AAGTACAATT CAAGAAACTG AGTATTTATG GGCATTGAAG AAAAAATGTT 1851 GAGATAAAAT TGCTGTGCAG AAAAAAGTGT TAATGAAGCC GACCTGACTA 10 1901 CTTAACCTTA GAGACCTGCT TTACAAGGTT GGCCCTTGAT TGGCATCTGG 1951 GAACTTGGAG TTCAGGGGGC TTCCACCATT CCCAGAACTG ATCAAAGTAG 2001 CTTACTATAT CTAAACTGTA AAACAATATA GTTTCTCCTG AACACCTGCT 2051 TTCCTTCTGG GAGTCTGGAA TTTTGGTATG TGCCAGGCAG AGACTACCTT PUSL TICCTICIGG GAGICIGGAA TITIGGTATG IGCCAGGCAG AGACTACCTI
PLOL TGTGACCAGC TCCCAGTAAA AACCCCAGGC ACTCAGTCTC TAACAAGCTT
PLOL TTCTGGTTGA CAGTGTTTCA CAAGTGCTGT TACAACTGGT TGCTGGGAGA
PROL TTCTGTGA TCCTCTGTGA TTCCACTGGC GGAGGATTCT TGGAAGCTTG
PROL TTTGTTTTGT ACCCTTGACTT CACCCCATGT GTCTTTTTTC CTTTGCTGAT
PROL TTTGTTTTGT ATCCTTTCAC TGTAATAAAA AAAAAA
PROL TTTGTTTTGT ATCCTTTCAC TGTAATAAAA AAAAAA 15 20

BLAST Results

25

No BLAST result

Medline entries

30

No Medline entry

35

Peptide information for frame 1

ORF from 46 bp to 702 bp; peptide length: 219
40 Category: similarity to unknown protein
Classification: unclassified

1 MNGTKALAIE EFCKSLGHAC IRFDYSGVGS SDGNSEESTL GKWRKDVLSI
51 IDDLADGPQI LVGSSLGGWL MLHAAIARPE KVVALIGVAT AADTLVTKFN
45 101 QLPVELKKEV EMKGVWSMPS KYSEEGVYNV QYSFIKEAEH HCLLHSPIPV
151 NCPIRLLHGM KDDIVPWHTS MQVADRVLST DVDVILRKHS DHRMREKADI
201 QLLVYTIDDL IDKLSTIVN

50

BLASTP hits

No BLASTP hits available

55 Alert BLASTP hits for DKFZphmel2_7kl9, frame l

No Alert BLASTP hits found

Pedant information for DKFZphmel2_7kl9, frame 1

Report for DKFZphmel2_7kl9.1

5

ELENGTHD 219

EMWD 24309.18

IpID 5-69

10 [HOMOL] PIR:A71691 hypothetical protein RP343 - Rickettsia

prowazekii 3e-29 EBLOCKSI BPO4352K EBLOCKSI PROD828E

EKWI Alpha_Beta

15

25

- - SEQ KYSEEGVYNVQYSFIKEAEHHCLLHSPIPVNCPIRLLHGMKDDIVPWHTSMQVADRVLST

 - SEQ DVDVILRKHSDHRMREKADIQLLVYTIDDLIDKLSTIVN PRD hheeeeecccchhhhhhhheeeeehhhhhhhhccccc
- 30 (No Prosite data available for DKFZphmel2_7kl9.1)
 - (No Pfam data available for DKFZphmel2_?kl9.l)
- 35 DKFZphtes3_10ilb

group: nucleic acid management

40

55

DKFZphtes3_10il6 encodes a novel 742 amino acid protein with similarity to human ZK1.

- The ZKL gene is one of early response genes by exposure to
 45 ionizing radiation, and plays a role in radiation-induced
 apoptotic cell death on hematopoietic cells. The novel protein
 contains 18 zinc finger domains, a RGD cell attachment and a ATP
 GTP A domain.
- The new protein can find application in diagnosis/therapy in leukemia predisposition/disease in the modulation of DNA repair.
 - similarity to ZKL (Homo sapiens), complete cds.

Sequenced by @iagen

Locus: unknown

-264-

Insert length: 2884 bp
Poly A stretch at pos. 2861, polyadenylation signal at pos. 2835

```
5
        1 CGGAAATGGA GGGGGTCGCT TTCCTCACCT TCCTCGCTGC GCGGGCGGCG
       51 GTTGGTAACC GGTCAGACCA GCCCGAGAGG GACCTGGTGC CTGTACCCAG
      LOL GCTTCTGTCG CTCTGTCGCC TGCGCTATGC CCTGCTGTAG TCACAGGAGC
      151 TGTAGAGAGG ACCCCGGTAC ATCTGAAAGC CGGGAAATGG ACCCAGTGGC
      201 CTTTGAGGAT GTGGCTGTGA ACTTCACCCA GGAAGAGTGG ACATTGCTGG
10
      251 ATATTTCCCA GAAGAATCTC TTCAGGGAAG TGATGCTGGA AACTTTCAGG
      ADI AACCTGACCT CTATAGGAAA AAAATGGAGT GACCAGAACA TTGAATATGA
      351 GTACCAAAAC CCCAGAAGAA GCTTCAGGAG TCTCATAGAA GAGAAAGTCA
      4D1 ATGAAATTAA AGAAGACAGT CATTGTGGAG AAACTTTTAC CCAGGTTCCA
      451 GATGACAGAC TGAACTTCCA GGAGAAGAAA GCTTCTCCTG AAGTAAAATC
15
      501 ATGTGACAGC TTTGTGTGTG CAGAAGTTGG CATAGGTAAC TCATCTTTTA
      55% ATATGAGCAT CAGAGGTGAC ACTGGACACA AGGCATATGA GTATCAGGAA
      LOL TATGGACCAA AGCCATATAA GTGTCAACAA CCTAAAAATA AGAAAGCCTT
      L51 CAGGTATCGC CCATCCATTA GAACACAAGA AAGGGATCAC ACTGGAGAGA
      701 AACCCTATGC TTGTAAAGTC TGTGGAAAAA CCTTTATTTT CCATTCAAGC
20
      751 ATTCGAAGAC ACATGGTAAT GCACAGTGGG GATGGAACTT ATAAATGTAA
      BD3 ATTTTGTGGG AAAGCCTTCC ATTCTTTCAG TTTATATCTT ATCCATGAAA
      B51 GAACTCACAC TGGAGAGAAA CCATATGAAT GTAAACAATG TGGTAAATCC
      POD TTTACTTATT CTGCTACCCT TCAAATACAT GAAAGAACTC ACACTGGGGA
      951 GAAGCCCTAT GAATGTAGCA AATGTGATAA AGCATTTCAT AGTTCTAGTT
     LODL CCTATCATAG ACATGAAAGA AGTCACATGG GAGAGAAGCC TTATCAATGC
     1051 AAAGAATGTG GAAAAGCATT TGCATATACC AGTTCTCTTC GTAGACATGA
     BIDI AAGGACCCAC TCTGGGAAAA AACCGTATGA ATGTAAGCAA TATGGGGAAG
     1151 GCTTATCCTA TCTTATAAGT TTTCAAACAC ACATAAGAAT GAACTCTGGA
     1201 GAAAGACCTT ATAAATGTAA GATATGTGGG AAAGGCTTTT ATTCTGCCAA
     1251 GTCATTTCAA ACACATGAAA AAACTCACAC TGGAGAGAAA CGCTATAAAT
     LIDI GCAAGCATG TGGTAAAGCC TTCAATCTTT CCAGTTCCTT TCGATATCAT
     1351 GAAAGGATTC ACACTGGAGA GAAACCCTAT GAGTGTAGC AGTGTGGGAA
     1401 AGCCTTCAGA TCTGCCTCAC AGCTTCGAGT GCACGGTGGG ACTCACACTG
     1451 GAGAGAAACC CTATGAATGT AAGGAATGTG GGAAAGCCTT CAGATCTACC
35
     1501 TCACACCTTC GAGTGCATGG TAGGACTCAT ACTGGAGAGA AACCCTATGA
     1551 ATGTAAGGAA TGTGGGAAAG CCTTCAGATA TGTGAAGCAC CTTCAAATTC
    JLOJ ATGAAAGGAC AGAAAAACAC ATAAGAATGC CCTCTGGAGA AAGACCTTAT
     1651 AAATGTAGTA TATGTGAGAA AGGCTTTTAT TCTGCCAAGT CATTTCAAAC
     1701 ACATGAAAAA ACTCACACTG GAGAGAAACC CTATGAATGC AACCAATGTG
40
    1751 GTAAAGCCTT CAGATGTTGC AATTCCCTTC GATATCATGA AAGGACTCAC
     1801 ACTGGAGAGA AACCCTATGA GTGTAAGCAA TGTGGGAAAG CCTTCAGATC
     1851 TGCCTCACAC CTTCGAATGC ATGAAAGGAC TCACACTGGA GAGAAACCCT
     1901 ATGAGTGTAA GCAATGTGGG AAAGCCTTCA GTTGTGCCTC AAACCTTCGA
1951 AAGCATGGTA GGACTCACAC TGGAGAGAAA CCCTATGAGT GTAAGCAATG
45
     2001 TGGGAAAGCC TTCAGATCTG CCTCAAACCT TCAGATGCAT GAAAGGACTC
     2051 ACACTGGAGA GAAACCCTAT GAATGTAAGG AATGCGAAAA AGCATTCTGT
     2101 AAATTCTCTT CTTTTCAAAT ACATGAAAGG AAGCACAGAG GAGAGAAGCC 2151 CTATGAATGT AAGCATTGTG GGAATGGATT CACATCTGCC AAGATTCTTC
     2201 AAATACATGC AAGAACACAC ATTGGAGAGA AACACTATGA ATGTAAGGAA
50
     2251 TGCGGAAAAG CATTCAATTA TTTTTCTTCC TTGCATATAC ACGCAAGGAC
     23Dl TCATATGGGA GAGAAGCCAT ATGAATGTAA GGATTGTGGG AAAGCATTCA
     2351 GCTAGCCTGG TTCCTTTTAT GGACATGAAT AGACTCACAC TGGAAGGAAG
     2401 CACTATGAAT GCAAGCAATG TGGCAAAACT TTCACATTTT CCAGTTCTTT
     2451 TCGATATCAT GAAAGGACTC ACACTGGGGA GAAACCCTAT CAATGTAAGC
55
     25DL AGTGTGGGAA AGCCTTCATT CCTTTTACTT CTTTTCAATG TCATGAAAGG
     2551 ACTCACACGG GAGAGAAACC CTATGAGTGT ATTCTAGTTC CGTTTGATAT
     2601 CATGAAAGGA CTTACACTGG AGTGAAACCC TATGAATGTA AGCAATGTGG
```

255 GAAAGCCTTC AGATGTGCCT CGCACCTTCA ACGGCATGGA AGGGTTCACA 2701 CTTGGGAGAA ACTCTATGAA TGTAAGCAGT ATGGGAAAGC CTTCAGATCT 2751 GCCAAGATTC TTTGAATACA GATAATTAAT GTAAACAATT ATCATAAGTA 2801 TACTAACATG TTATTCTTTT TAAATAAGAA GGTATAATAA AATATCCCAT 5 2851 TGGTTTTATG TATTAAAAAA AAAAAAAAA AAAA

BLAST Results

10

No BLAST result

Medline entries

15

98401134:

Katoh 0, Oguri T, Takahashi T, Takai S, Fujiwara Y, Watanabe H.; ZKl, a

- 20 novel Kruppel-type zinc finger gene, is induced following exposure to ionizing radiation and enhances apoptotic cell death on hematopoietic cells. Biochem Biophys Res Commun 1998 Aug 28:249(3):595-600
- 25 95137393:
 Wick MJ, Ann DK, Lee NM, Loh HH.; Isolation of a cDNA encoding a novel
 zinc-finger protein from
 neuroblastoma x glioma NG108-15 cells. Gene 1995 Jan
 30 23;152(2):227-32

Peptide information for frame 1

ORF from 127 bp to 2352 bp; peptide length: 742
Category: similarity to known protein

40 Classification: Nucleic acid management
Prosite motifs: RGD (146-148)
ATP_GTP_A (195-202)
ZINC_FINGER_C2H2 (196-216)
ZINC_FINGER_C2H2 (224-244)

45 ZINC_FINGER_C2H2 (252-272)
ZINC_FINGER_C2H2 (280-300)

ZINC_FINGER_C2H2 (354-384) ZINC_FINGER_C2H2 (372-412) 50 ZINC_FINGER_C2H2 (420-440)

ZINC_FINGER_C2H2 (3D8-328)

ZINC_FINGER_C2H2 (448-468)
ZINC_FINGER_C2H2 (510-530)
ZINC_FINGER_C2H2 (538-558)

ZINC_FINGER_C2H2 (566-586)

55 ZINC_FINGER_C2H2 (594-614)

ZINC_FINGER_C2H2 (622-642)

ZINC_FINGER_C2H2 (650-670)

ZINC_FINGER_C2H2 (678-698)

ZINC_FINGER_C2H2 (706-726) ZINC_FINGER_C2H2 (476-498)

5 1 MPCCSHRSCR EDPGTSESRE MDPVAFEDVA VNFTQEEUTL LDISQKNLFR 51 EVMLETFRNL TSIGKKWSDQ NIEYEYQNPR RSFRSLIEEK VNEIKEDSHC 101 GETFTQVPDD RLNFQEKKAS PEVKSCDSFV CAEVGIGNSS FNMSIRGDTG 151 HKAYEYGEYG PKPYKCQQPK NKKAFRYRPS IRTQERDHTG EKPYACKVCG 201 KTFIFHZZIR RHMVMHZGDG TYKCKFCGKA FHZFZLYLIH ERTHTGEKPY 251 ECKQCGKSFT YSATLQIHER THTGEKPYEC SKCDKAFHSS SSYHRHERSH 10 3D1 MGEKPYQCKE CGKAFAYTSS LRRHERTHSG KKPYECKQYG EGLSYLISFQ 351 THIRMNSGER PYKCKICGKG FYSAKSFRTH EKTHTGEKRY KCKRCGKAFN 401 LSSSFRYHER IHTGEKPYEC KQCGKAFRSA SQLRVHGGTH TGEKPYECKE 451 CGKAFRSTSH LRVHGRTHTG EKPYECKECG KAFRYVKHL@ IHERTEKHIR 501 MPSGERPYKC SICEKGFYSA KSFQTHEKTH TGEKPYECNQ CGKAFRCCNS 15 551 LRYHERTHTG EKPYECKQCG KAFRSASHLR MHERTHTGEK PYECKQCGKA LOT LECAZNERKH GRIHTGEKPY ECKQCGKAFR ZAZNEQMHER THIGEKPYEC L51 KECEKAFCKF SSF@IHERKH RGEKPYECKH CGNGFTSAKI L@IHARTHIG 701 EKHYECKECG KAFNYFSSLH IHARTHMGEK PYECKDCGKA FS

20

BLASTP hits

25 No BLASTP hits available

Alert BLASTP hits for DKFZphtes3_10ilb, frame 1

No Alert BLASTP hits found

30.

Peptide information for frame 2

35 ORF from 1703 bp to 2584 bp; peptide length: 294 Category: questionable ORF Classification: no clue

1 MKKLTLERNP MNATNVVKPS DVAIPFDIMK GLTLERNPMS VSNVGKPSDL
51 PHTFECMKGL TLERNPMSVS NVGKPSVVPQ TFESMVGLTL ERNPMSVSNV
101 GKPSDLPQTF RCMKGLTLER NPMNVRNAKK HSVNSLLFKY MKGSTEERSP
151 MNVSIVGMDS HLPRFFKYMQ EHTLERNTMN VRNAEKHSII FLPCIYTQGL
201 IWERSHMNVR IVGKHSASLV PFMDMNRLTL EGSTMNASNV AKLSHFPVLF
251 DIMKGLTLGR NPINVSSVGK PSFLLLLFNV MKGLTRERNP MSVF

45

BLASTP hits

50 No BLASTP hits available

Alert BLASTP hits for DKFZphtes3_10ilb, frame 2

TREMBL:AF153201_1 product: "zinc finger protein dp"; Homo
55 sapiens zinc
finger protein dp mRNA, complete cds, N = 1, Score = 225, p =
4.le-l8

>TREMBL:AF153201_1 product: "zinc finger protein dp"; Homo sapiens zinc finger protein dp mRNA, complete cds. Length = 423 5 HSPs: Score = 225 (33.8 bits), Expect = 4.le-l8, P = 4.le-l8 Identities = 84/246 (34%), Positives = 122/246 (49%) 10 Query: 16 VVKPSDVA-IPFDIMKGLTLERNPMSVSNVGKPSDLPHTFECMKGLTLERNPMSVSNVGK 74 V KPS A I F I + + L RN + V +V K S TLERNP++V +VGK - IZBRIJAJUA SVENJEV E LGRNHIHVISVAKVSVRIQTLLNIEGSTLERNPINVMSVGK 61 PSVVPQTFESMVGLTLERNPMSVSNVGKPSDLPQTFRCMKGLTLERNPMVRAKKHSVN 134 20 + Q+ + G LERNP+ V NV KPS Sbjct: LLIRAQSLFYIRGFILERNPIPVINVAKPSVGFQILLIINEFTLERSLTHVISAIKCLVE 121 25 135 SLLFKYMKGSTEERSPMNVSIVGMDS-HLPRFFKYM@EHTLERNTMNVRNAEKHSIIFLP 193 + R+PMNV VG P F +++E TLERN M+V 122 DEILLNITEFIQVRNPMNVMNVGKPLVRAPTLF-30 Sbjct: FIRESTLERNLMHVVIVLKALVAVQI 180 Querv: CIYTQGLIWERSHMNVRIVGKHSASLVPFMDMNRLTLEGSTMNASNVAKLSHFPVLFDIM 253 35 ER+HM+V V K +++ TL S + A V K S Sbjct: 181 LLSIKEYTLERNHMHVISVIKVLVKAQTSLNIREYTLVKSLIIAIVVRKPSVRVLTLFFI 240 254 KGLTLGRN 261 40 Query: + TL +N 241 REFTLEKN 248 Sbjct: Score = 215 (32.3 bits), Expect = 1.le-lb, P = 1.le-lb45 Identities = 82/246 (33%), Positives = 124/246 (50%) Query: VGKPSDLPHTFECMKGLTLERNPMSVSNVGKPSVVPQTFESMVGLTLERNPMSVSNVGKP 103 C++ L RN + V +V K SV QT ++ G VGKPS 50 TLERNP++V +VGK 3 Sbict: VGKPSVRAQILFCIRESILGRNHIHVISVAKVSVRIQTLLNIEGSTLERNPINVMSVGKL 62 Querv: 104 SDLPQTFRCMKGLTLERNPMNVRNAKKHSVNSLLFKYMKGSTEERSPMNV-55 SIVGM---D 159

-268-

++G LERNP+ V N K SV +

T ERS +V S

Q+

D

Sbjct: 63

LIRAQSLFYIRGFILERNPIPVINVAKPSVGFQILLIINEFTLERSLTHVISAIKCLVED 122

Query: 160 SHLPRFFKYMQEHTLERNTMNVRNAEKHSIIFLPCIY-

5 TQGLIWERSHMNVRIVGKHSAS 218

L +++Q RN MNV N K ++ P ++ + ER+ M+V

Sbjct: 123 EILLNITEFIQV----RNPMNVMNVGK-PLVRAPTLFFIRESTLERNLMHVVIVLKALVA 177

10

Query: 219 LVPFMDMNRLTLEGSTMNASNVAK-LSHFPVLFDIMKGLTLGRNPINVSSVGKPSFLLLL 277

+ + + TLE + M+ +V K L +I + TL ++ I V

KPS +L

15 Sbjct: 178 V@ILLSIKEYTLERNHMHVISVIKVLVKA@TSLNIRE-YTLVKSLIIAIVVKAPSVRVLT 236

Query: 278 FNVMKGLTRERN 289

++ T E+N

20 Sbjct: 237 LFFIREFTLEKN 248

Score = 207 (31.1 bits), Expect = 5.2e-15, P = 5.2e-15 Identities = 80/270 (29%), Positives = 129/270 (47%)

25 Query: 1 MKKLTLERNPMNATNVVKPSDVAIPFDIMKGLTLERNPMSVSNVGKPSDLPHTFECMKG 59
+++ L RN ++ ++G TLERNP++V +VGK

Sbjct: Lb IRESILGRNHIHVISVAKVS-

30 VRIQTLLNIEGSTLERNPINVMSVGKLLIRAQSLFYIRG 74

Query: 60

LTLERNPMSVSNVGKPSVVPQTFESMVGLTLERNPMSVSNVGKPSDLPQTFRCMKGLTLE .139

LERNP+ V NV KPSV Q + TLER+ V + K +

35

Sbjct: 75
FILERNPIPVINVAKPSVGFQILLIINEFTLERSLTHVISAIKCLVEDEILLNITEFIQV 134

Query: 120

40 RNPMNVRNAKKHSVNSLLFKYMKGSTEERSPMNVSIVGMDSHLPRFFKYMQEHTLERNTM 179
RNPMNV N K V + +++ ST ER+ M+V IV +

++E+TLERN M

Sbjct: 135

RNPMNVMNVGKPLVRAPTLFFIRESTLERNLMHVVIVLKALVAVQILLSIKEYTLERNHM 194

45

Query: 180
NVRNAEKHSIIFLPCIYTQGLIWERSHMNVRIVGKHSASLVPFMDMNRLTLEGSTMNASN 239
+V + K + + + + + C + + C + + C + + C + + C + + C + + C + + C + + C + + C + + C + + C + + C + + C + + C + + C + + C

50 Sbjct: 195

HVISVIKVLVKAQTSLNIREYTLVKSLIIAIVVRKPSVRVLTLFFIREFTLEKNYYLCTQ 254

Query: 240 VAKLSHFPVLFDIMKGLTL--GRNPINVSSVGK 270

+K F + D++K + G P S K

55 Sbjct: 255 CSK--SFSQISDLIKHQRIHTGEKPYKCSECRK 285

Score = l81 (27.2 bits) = Expect = l.4e-ll = P = l.4e-ll
Identities = 74/269 (27%) = Positives = ll6/269 (43%)

Query: TLERNPMNATNVVKPSDVAIPFDIMKGLTLERNPMSVSNVGKPSDLPHTFECMKGLTLER 64 TLERNP+N +V K A ++G LERNP+ V NV KPS TLER Sbjct: 48 TLERNPINVMSVGKLLIRAQSLFYIRGFILERNPIPVINVAKPSVGFQILLIINEFTLER 107 NPMSVSNVGKPSVVPQTFESMVGLTLERNPMSVSNVGKPSDLPQTFRCMKGLTLERNPMN 124 10 V + K V + ++ RNPM+V NVGKP TLERN M+ Sbjct: 108 SLTHVISAIKCLVEDEILLNITEFIQVRNPMNVMNVGKPLVRAPTLFFIRESTLERNLMH 167 15 125 VRNAKKHSVNSLLFKYMKGSTEERSPMNV-SIVGMDSHLPRFFKYMQEHTLERNTMNVRN 183 κv + +K T ER+ M+V S++ + ++E+TL ++ + 168 VVIVLKALVAVQILLSIKEYTLERNHMHVISVIKVLVKAQTSLN-Sbjct: 20 IREYTLVKSLIIAIV 226 184 Query: AEKHSIIFLPCIYTQGLIWERSHWNVRIVGKHSASUNDTRUTTEGSTWASNVAKL 243 25 K Z+ L + + E+++ K + + Sbjct: 227 VRKPSVRVLTLFFIREFTLEKNYYLCTQCSKSFSQISDLIKHQRIHTGEKPYKCSECRKA 286 244 SHFPVLFDIMKGLTLGRNPINVSSVGKPSF 273 30 Querv: L + + + G+ P GK SF 287 FSQCSLLALHQRIHTGKKPNPCDECGK-SF 315 Sbjct: Score = 166 (24.9 bits), Expect = 8.4e-10, P = 8.4e-10 35 Identities = 63/194 (32%), Positives = 89/194 (45%) Query: VGKPSDLPQTFRCMKGLTLERNPMNVRNAKKHSVNSLLFKYMKGSTEERSPMNVSIVGMD 159 VGKPS Q C++ L RN ++V + K SV 40 ER+P+NV VG Sbict: VGKPSVRAQILFCIRESILGRNHIHVISVAKVSVRIQTLLNIEGSTLERNPINVMSVGKL 62 Query: 760 SHLPRFFKYMQEHTLERNTMNVRNAEKHSIIFLPCIYTQGLIWERSHMNVRIVGKHSASL 219 45 Y++ LERN + V N K S+ F K Sbjct: 63 LIRAQSLFYIRGFILERNPIPVINVAKPSVGFQILLIINEFTLERSLTHVISAIKCLVED 122 50 220 VPFMDMNRLTLEGSTMNASNVAK-Query: LSHFPVLFDIMKGLTLGRNPINVSSVGKPSFLLLLF 278 + MN NV K L P LF I + TL RN ++V V K 123 EILLNITEFIQVRNPMNVMNVGKPLVRAPTLFFIRES~ 55 Sbjct: TLERNLMHVVIVLKALVAVQIL 181

Query: 279 NVMKGLTRERNPMSV 293

+K T ERN M V Sbjct: 182 LSIKEYTLERNHHHV 196

5 Pedant information for DKFZphtes3_10ilb, frame 1

Report for DKFZphtes3_10il6.1

10 **ELENGTHI** 784 EMWI 90857-05 [[pI] 9.24 EHOMOLI TREMBL: ABOLL414_1 gene: "ZKl"; product: "Kruppeltype zinc finger protein"; Homo sapiens ZKL mRNA for Kruppel-type zinc finger protein, complete cds. []. **IFUNCATI** 30-10 nuclear organization ES. cerevisiae, YJLD5bcl EE-93 **EFUNCATI** 04.05.01.04 transcriptional control ES. cerevisiae. 20 YJLO56cJ 6e-33 **EFUNCATI** 04.99 other transcription activities [S. cerevisiae] YOR113w1 5e-24 EFUNCATI 04.01.01 rrna synthesis ES. cerevisiae, YPR186c PZF1 -TFIIIAD le-20 EFUNCATI 04.03.01 trna synthesis ES. cerevisiae, YPRlabc PZFL -TFIIIAD le-20 **EFUNCATI** 13.04 homeostasis of other ions [S. cerevisiae, EL-91 Ew7501NY EFUNCATI 11.07 detoxification IS. cerevisiae, YGL254w1 2e-12 [FUNCAT] 01-02-04 regulation of nitrogen and sulphur utilization 30 ES. cerevisiae, YGL254wl 2e-l2 **EFUNCATI** 01.05.04 regulation of carbohydrate utilization cerevisiae, YGL209w1 2e-11 **EFUNCATI** 04.05.99 other mrna-transcription activities EZ-35 cerevisiae, YERO28cl 3e-10 [FUNCAT] 11-01 stress response ES. cerevisiae, YKLOb2wl le-09 EFUNCATE 01.01.04 regulation of amino-acid metabolism cerevisiae, YDR253cJ 5e-D9 **EFUNCATE** 99 unclassified proteins ES. cerevisiae, YBRD66cl 40 80-9E EFUNCATI D3.07 pheromone response, mating-type determination, **EFUNCATI** BLOD466 TFIIS zinc ribbon domain proteins **EBFOCK21** 45 BLOD245A Phytochrome chromophore attachment site **EBFOCK21** proteins EBF0CK21 DM019518 EBF0CK23 PF01363B EBF0CK23 BF07030 **CBLOCK21** PF00096B 50 BLODD28 Zinc finger, C2H2 type, domain proteins **EBFOCK23 EBFOCK2**3 BP04213E **EBFOCK2** BP04233C EBFOCKZI Bb04573B 55 d2adr__ 7.31.1.4 ADR1 Esynthetic based on yeast **EZCOPI** (Saccharomyce 2e-05 nucleus le-53 CPIRKW] **EPIRKU**I RNA binding 2e-58

WO 01/98454 PCT/IB01/02050 **EPIRKWI** duplication le-34 **EPIRKU**3 tandem repeat le-171 **EPIRKUJ** spermatogenesis 5e-62 zinc le-169 **EPIRKUJ** zinc finger 0.0 **CPIRKWI EPIRKUI** DNA binding D.O metal binding le-120 **EPIRKWI EPIRKUJ** phosphoprotein 2e-58 leucine zipper le-53 **EPIRKWI** alternative splicing 2e-58 10 **EPIRKU** EPIRKWI EPIRKWI eye lens le-lll oocyte le-106 transcription factor le-lll **EPIRKU** embryo le-106 [PIRKW] segmentation le-34 15 **EPIRKUJ** [PIRKU] transcription regulation le-152 ESUPFAMD POZ domain homology 7e-83 ESUPFAMI transcription factor Krueppel le-34
ESUPFAMI zinc finger protein ZFP-36 le-173 **ESUPFAM1** transcription factor IIIA 8e-31 20 IPROSITED ATP_GTP_A 1 **EPROSITED** RGD EPROSITED ZINC_FINGER_C2H2 Zinc finger, C2H2 type **EPFAMI** 25 [PFAM] TNFR/NGFR cysteine-rich region EKWI Irregular EKWI 3 D LOW_COMPLEXITY EKWI 3.57 % 30 RKWRGSLSSPSSLRGRRLVTGQTSPRGTWCLYPGFCRSVACAMPCCSHRSCREDPGTSES SEQ SEG **l**meyF 35 REMDPVAFEDVAVNFTQEEWTLLDISQKNLFREVMLETFRNLTSIGKKWSDQNIEYEYQN SEQ SEG lmeyF 40 PRRSFRSLIEEKVNEIKEDSHCGETFTQVPDDRLNFQEKKASPEVKSCDSFVCAEVGIGN SEQ SEG lmeyF 45 **SSFNMSIRGDTGHKAYEY@EYGPKPYKC@@PKNKKAFRYRPSIRT@ERDHTGEKPYACKV** SEQ SEG lmeyF 50 CGKTFIFHSSIRRHMVMHSGDGTYKCKFCGKAFHSFSLYLIHERTHTGEKPYECKQCGKS SEG lmevF 55 SEQ FTYSATLQIHERTHTGEKPYECSKCDKAFHSSSSYHRHERSHMGEKPY@CKECGKAFAYT

SEG

1			
			÷

	lmeyF					· · · · · · · · · · · · · · · · · · ·	
5	SEQ SEG lmeyF		• • • • • • • • • •	• • • • • • • • •	FQTHIRMNSGER	• • • • • • • • • • • •	
10	SEQ SEG lmeyF	• • • • • • •	• • • • • • • • • •	• • • • • • • • • • •	ERIHTGEKPYECI	• • • • • • • • • • • • •	• • • • •
15	SEQ SEG ImeyF	• • • • • • •			TGEKPYECKECG	• • • • • • • • • • • • • • • • • • • •	• • • • •
20	SEQ SEG lmeyF	•••••		• • • • • • • • • • •	THTGEKPYECNQ	• • • • • • • • • • • •	• • • • •
25	SEQ SEG lmeyF	• • • • • • • •		• • • • • • • • • •	TCCEEETTTTEE	• • • • • • • • • • • • •	• • • • •
30	SEG lmeyF	• • • • • • • •	• • • • • • • • • •	• • • • • • • • • • • •	ECKECEKAFCKF	• • • • • • • • • • • • • • • • • • • •	• • • • •
35	SEQ SEG lmeyF	• • • • • • •	• • • • • • • • • • •	• • • • • • • • • •	CGKAFNYFZSLH	• • • • • • • • • • • • • • • • • • • •	• • • • • •
40		KAFS ····	. *				
45			Pro	site for DK	:FZphtes3_10i	16.1	
50	PS000 PS000 PS000 PS000 PS000 PS000 PS000 PS000	17 28 28 28 28 28 28	188->191 237->245 238->259 266->287 294->315 322->343 350->371 406->427 434->455	RGD ATP_GTP_A ZINC_FINGE ZINC_FINGE ZINC_FINGE ZINC_FINGE ZINC_FINGE ZINC_FINGE ZINC_FINGE	:R_C2H2 :R_C2H2 :R_C2H2 :R_C2H2	PD0COOO36 PD0COOO38 PD0COOO38 PD0COOO38 PD0COOO38 PD0COOO38 PD0COOO38	
,,	P2000 P2000 P2000	28 28	454-2433 452-243 490-2511 552-2573	ZINC_FINGE ZINC_FINGE ZINC_FINGE	:K_C5H5 :K_C5H5	PD0C00028 PD0C00028 PD0C00028	

```
WO 01/98454
                                                           PCT/IB01/02050
    - 850002A
                               ZINC_FINGER_C2H2
                   580->601
                                                          PD0C00058
                               ZINC_FINGER_C2H2
                   608->629
    $500024
                                                          PD0000028
                   636->657
                               ZINC_FINGER_C2H2
    820002A
                                                          PD0C00058
                               ZINC_FINGER_C2H2
                   664->685
    B200028
                                                          PD0C00028
                               ZINC_FINGER_C2H2
    B200059
                   692->713
                                                          PD0CDDD28
                               ZINC_FINGER_C2H2
                   720->743
    9200024
                                                          PD0C00058
                   748->769
                               ZINC_FINGER_C2H2
    850002A
                                                          PD0C00058
                               ZINC_FINGER_C2H2
    9200024
                   518->541
                                                          PD0C00058
10
                            Pfam for DKFZphtes3_10il6.1
15
    HMM_NAME TNFR/NGFR cysteine-rich region
    HMM
                         *CpeGtYtD.WNHvpqClpC..trCePEMGQYMvqPCTwTQNTVC*
                          C + +++ +++++C C ++C+++ G++++++ ++
                         CLYPGFCRSVACAMPC--CSHRSCREDPGTSESREMDP----VA
    Query
20
    67
    HMM_NAME Zinc finger, C2H2 type
25
    MMH
                         *CpwPDCaKtFrrwsNLrRHMRTH*
                          C++ CGKTF
                                      S+ RRHM +H
    Query
                          CKV--CGKTFIFHSSIRRHMVMH
                                                       258
    32.15 (bits) f: 266 t: 286 Target: dkfzphtes3_10il6.l similarity to ZK1 (Homo sapiens), complete cds.
30
      Alignment to HMM consensus:
    Query
                         *CpwPDCgKtFrrwsNLrRHMRTH*
                          C++ CGK+F + S + +H RTH
                    266 CKF--CGKAFHSFSLYLIHERTH
35
      dkfzphtes3
                                                       286
                   f: 294 t: 314 Target: dkfzphtes3_10i16.1
    similarity to ZK1 (Homo sapiens), complete cds.
      Alignment to HMM consensus:
40
                         *CpwPDCgKtFrrwsNLrRHMRTH*
                              CGK+F+++ +L++H RTH
                          C+
                    294 CKQ--CGKSFTYSATLQIHERTH
    Query
                                                       314
    34.22 (bits) f: 322 t: 342 Target: dkfzphtes3_10i16.1
    similarity to ZKL (Homo sapiens), complete cds.
45
      Alignment to HMM consensus:
                         *CpwPDCgKtFrrwsNLrRHMRTH*
    Query
                          C++ C+K+F ++S++ RH R+H
                    322 CSK--CDKAFHSSSSYHRHERSH
      dkfzphtes3
                                                       342
50
    Query f: 350 t: 370 Target: dkfzphtes3_l0il6.l
similarity to ZKl (Homo sapiens), complete cds.
      Alignment to HMM consensus:
                         *CpwPDCgKtFrrwsNLrRHMRTH*
    HMM
55
                          C++ CGK+F + S+LRRH RTH
```

370

350 CKE--CGKAFAYTSSLRRHERTH

Querv

PCT/IB01/02050 WO 01/98454

32.09 (bits) f: 406 t: 426 Target: dkfzphtes3_10i16.1 similarity to ZKL (Homo sapiens), complete cds. Alignment to HMM consensus: *CpwPDCgKtFrrwsNLrRHMRTH* Query C++ CGK F ++ ++++H +TH 5 CKI--CGKGFYSAKSFQTHEKTH dkfzphtes3 406 426 f: 434 t: 454 Target: dkfzphtes3_10ilb.l Query similarity to ZKL (Homo sapiens), complete cds. 10 Alignment to HMM consensus: *CpwPDCqKtFrrwsNLrRHMRTH* C+ CGK+F+ +S++R H R+H Query 434 CKQ--CGKAFNLSSSFRYHERIH 454 32.94 (bits) f: 462 t: 482 Target: dkfzphtes3_10i16.1 15 similarity to ZKL (Homo sapiens), complete cds. Alignment to HMM consensus: *CpwPDCgKtFrrwsNLrRHMRTH* Query CGK+FR++S+LR H TH (+ 20 dkfzphtes3 462 CKQ--CGKAFRSASQLRVHGGTH 482 f: 490 t: 510 Target: dkfzphtes3_10i16.1 **Query** similarity to ZKL (Homo sapiens), complete cds. Alignment to HMM consensus: *CpwPDCgKtFrrwsNLrRHMRTH* 25 C++ CGK+FR+ S+LR H RTH 490 CKE--CGKAFRSTSHLRVHGRTH 510 Query 30.69 (bits) f: 518 t: 540 Target: dkfzphtes3_10i16.1 similarity to ZK1 (Homo sapiens), complete cds. 30 Alignment to HMM consensus: *CpwPDCgKtFrrwsNLrRHMR..T.H* C++ CGK+FR+ +L++H R 518 CKE--CGKAFRYVKHLQIHERTE-KH 540 dkfzphtes3 35 Query f: 552 t: 572 Target: dkfzphtes3_10i16.1 similarity to ZK1 (Homo sapiens), complete cds. Alignment to HMM consensus: *CpwPDCqKtFrrwsNLrRHMRTH* HMM C++ C+K F ++ ++++H +TH 40 552 CSI--CEKGFYSAKSFQTHEKTH 572 Query 31.33 (bits) f: 580 t: 600 Target: dkfzphtes3_10il6.1 similarity to ZK1 (Homo sapiens), complete cds. 45 Alignment to HMM consensus: *CpwPDCgKtFrrwsNLrRHMRTH* Query C+ CGK+FR +LR H RTH dkfzphtes3 580 CNQ--CGKAFRCCNSLRYHERTH P00

50 f: 608 t: 628 Target: dkfzphtes3_10i16.1 Querv similarity to ZK1 (Homo sapiens), complete cds. Alignment to HMM consensus: *CpwPDCqKtFrrwsNLrRHMRTH* MMH CGK+FR++S+LR+H RTH **C**+

CKQ--CGKAFRSASHLRMHERTH 55 604 P58 Query

35.30 (bits) f: 636 t: 656 Target: dkfzphtes3_10il6.1 similarity to ZK1 (Homo sapiens), complete cds.

WO 01/98454 PCT/IB01/02050 Alignment to HMM consensus: Query *CpwPDCgKtFrrwsNLrRHMRTH* C+ CGK+F+ +SNLR+H RTH dkfzphtes3 636 CKQ--CGKAFSCASNLRKHGRTH **656** 5 f: 664 t: 684 Target: dkfzphtes3_10i16.1 Query similarity to ZK1 (Homo sapiens), complete cds. Alignment to HMM consensus: MMH *CpwPDCgKtFrrwsNLrRHMRTH* 10 C+ CGK+FR++SNL++H RTH Querv 664 CKQ--CGKAFRSASNLQMHERTH 684 31.74 (bits) f: 692 t: 712 Target: dkfzphtes3_10i16.1 similarity to ZKl (Homo sapiens), complete cds. 15 Alignment to HMM consensus: Query *CpwPDCgKtFrrwsNLrRHMRTH* C++ C+K+F+ S+++H R H dkfzphtes3 692 CKE--CEKAFCKFSSF@IHERKH 712 20 Query 720 t: 740 Target: dkfzphtes3_10i16.1 f: similarity to ZKL (Homo sapiens), complete cds. Alignment to HMM consensus: *CpwPDCqKtFrrwsNLrRHMRTH* C++ CG F+++ L++H RTH 25 720 CKH--CGNGFTSAKILQIHARTH Query 740 748 t: 768 Target: dkfzphtes3_10i16.1 34.88 (bits) f: similarity to ZKL (Homo sapiens), complete cds. Alignment to HMM consensus: 30 Query *CpwPDCgKtFrrwsNLrRHMRTH* C++ CGK+F++ S+L +H RTH dkfzphtes3 748 CKE--CGKAFNYFSSLHIHARTH **768** 35 Pedant information for DKFZphtes3_10ilb, frame 2 Report for DKFZphtes3_10i16.2 40 **ELENGTHI** 294 CMMJ 33083.98 [[q] 9.97 45 EHOMOLI TREMBL:AF153201_1 product: "zinc finger protein dp": Homo sapiens zinc finger protein dp mRNA: complete cds. 7e-17 EKWI All_Alpha 50 MKKLTLERNPMNATNVVKPSDVAIPFDIMKGLTLERNPMSVSNVGKPSDLPHTFECMKGL SEQ ccccccccccccchhhhhcccccccccccccccchhhhhee PRD SEQ TLERNPMSVSNVGKPSVVPQTFESMVGLTLERNPMSVSNVGKPSDLPQTFRCMKGLTLER 55 PRD

NPMNVRNAKKHSVNSLLFKYMKGSTEERSPMNVSIVGMDSHLPRFFKYMQEHTLERNTMN

SEQ

PRD

	SEQ PRD	VRNAEKHSIIFLPCIYT&GLIWERSHMNVRIVGKHSASLVPFMDMNRLTLEGSTMNASNV chhhhhhheeeccceeeechhhhhhhcccccccc
5	SEQ PRD	AKLSHFPVLFDIMKGLTLGRNPINVSSVGKPSFLLLLFNVMKGLTRERNPMSVF ccccccchhhhhhhhccccccccccccchhhhhhhhhcccc
10	(No	Prosite data available for DKFZphtes3_10i16.2)
10	(No	Pfam data available for DKFZphtes3_10ilb.2)

DKFZphtes3_10n10

5 group: testis derived

DKFZphtes3_10nlO encodes a novel 502 amino acid protein without similarity to known proteins.

- 10 The mRNA is differentially polyadenylated and the novel protein is ubiquitously expressed. No informative BLAST results; No predictive prosite, pfam or SCOP motife.
- 15 The new protein can find application in studying the expression profile of testis-specific genes.

unknown protein

differentially polyadenylated

Sequenced by @iagen

25 Locus: unknown

20

Insert length: 255% bp
Poly A stretch at pos. 253%, polyadenylation signal at pos. 25%3

30 L CTCAGCCTCC CAAGTGGCTG GGACTGCAGG TTCTAAATGG CTTCTAAGAA 51 GTTGGGTGCA GATTTTCATG GGACTTTCAG TTACCTTGAT GATGTCCCAT LOL TTAAGACAGG AGACAAATTC AAAACACCAG CTAAAGTTGG TCTACCTATT LSL GGCTTCTCCT TGCCTGATTG TTTGCAGGTT GTCAGAGAAG TACAGTATGA 201 CTTCTCTTTG GAAAAGAAAA CCATTGAGTG GGCTGAAGAG ATTAAGAAAA 35 251 TCGAAGAAGC CGAGCGGGAA GCAGAGTGCA AAATTGCGGA AGCAGAAGCT DTDTTDDADT AAAADAAAA TGAGAGGG CCCAGAGGGC GATAGAAAA TGAGCTTCTC 351 CAAGACTCAC AGTACAGCCA CAATGCCAC TCCTATAAC CCCATCCTCG 401 CCAGCTTGCA GCACAACAGC ATCCTCACAC CAACTCGGGT CAGCAGTAGT
451 GCCACGAAAC AGAAAGTTCT CAGCCCACCT CACATAAAGG CGGATTTCAA
501 TCTTGCTGAC TTTGAGTGTG AAGAAGACCC ATTTGATAAT CTGGAGTTAA 40 551 AAACTATTGA TGAGAAGGAA GAGCTGAGAA ATATTCTGGT AGGAACCACT LOT GGACCCATTA TGGCTCAGTT ATTGGACAAT AACTTGCCCA GGGGAGGCTC LSL TGGGTCTGTG TTACAGGATG AGGAGGTCCT GGCATCCTTG GAACGGGCAA 45 751 CAGTTGGGCA ACTGTGAAAA GATGTCACTG TCTTCCAAAG TGTCCCTCCC BOD CCCTATACCT GCAGTAAGCA ATATCAAATC CCTGTCTTTC CCCAAACTTG 851 ACTCTGATGA CAGCAATCAG AAGACAGCCA AGCTGGCGAG CACTTTCCAT 901 AGCACATCCT GCCTCCGCAA TGGCACGTTC CAGAATTCCC TAAAGCCTTC 951 CACCCAAAGC AGTGCCAGTG AGCTCAATGG GCATCACACT CTTGGGCTTT 50 LOOL CAGCTTTGAA CTTGGACAGT GGCACAGAGA TGCCAGCCCT GACATCCTCC 1051 CAGATGCCTT CCCTCTCTGT TTTGTCTGTG TGCACAGAGG AATCATCACC ILDL TCCAAATACT GGTCCCACGG TCACCCCTCC TAATTTCTCA GTGTCACAAG 1151 TGCCCAACAT GCCCAGCTGT CCCCAGGCCT ATTCTGAACT GCAGATGCTG 1201 TCCCCCAGCG AGCGGCAGTG TGTGGAGACG GTGGTCAACA TGGGCTACTC 55 1251 GTACGAGTGT GTCCTCAGAG CCATGAAGAA GAAAGGAGAG AATATTGAGC 1301 AGATTCTCGA CTATCTCTTT GCACATGGAC AGCTTTGTGA GAAGGGCTTC LBSL GACCCTCTTT TAGTGGAAGA GGCTCTGGAA ATGCACCAGT GTTCAGAAGA

WO 01/98454 PCT/IB01/02050 1401 AAAGATGATG GAGTTTCTTC AGTTAATGAG CAAATTTAAG GAGATGGGCT BUSE TTGAGCTGAA AGACATTAAG GAAGTTTTGC TATTACACAA CAATGACCAG 15D1 GACAATGCTT TGGAAGACCT CATGGCTCGG GCAGGAGCCA GCTGAGACCA 1551 GGCCCTGCCT AGGCCCTGCC GCAGAACCAC CATCCCTGGG AGGCCCTGCA 1601 GAGCCCACCT GTGGGGAAAG AGAAGGGGCA GCTTCCGGAT TTTCTTTTGG 5 1651 GGGTTAGAAG GTCAGGTGTG GAGACTGCTC GCCAGTCTCT GTGAGCCTAG 17D1 GCCCTGAGCT GGGGAGGTGG GGAAGATTCG GGCATGTGAG TGCCCCCAGA 1751 ACTGTCCTGG CTCCTTCCGT ATTAAACGCA TTTGCATTTT GAGAAGTGTC 18D1 CTTCCCACTT CAGCCCTCCG GAGAGACTAC CCTAGTCTTT CTGGGGTGTT 1851 TATGTCCTCA GCTGAAGCCT GGCCTAGTTG CTGAGAGGGG CTGGGGAGAT 10 1901 GGGGCGGGAG GGCCAGACTC AGTGCTGCTG TGGAGCTAGG TGCTTCCCCC 1951 TTCCCCTGAG ACTGGTTGAC TGAACTCCAG TCAAGTTGAG TTCAAGTGAA 2001 AGATTCTTCC AGGGTTTTAT TTTTTCCCCT CCTAACAAAG TCTCATAGTG 2051 TTAACACTGG TTCTGCAATA TCTCTGAGGT GCAAAGAATG CACTTTTCCC 2101 TATGGGGCCC AGAGTTTGCC TTTTCTGCCA GGCAGTCACC ACGCTTCCCT 15 2151 ACCCCAGCCT GTTTCTTTTG GCTTGGTTTG GACCACAGTC CTCTGCTACC 2201 CAGGGTTTTA GAGCCCCTGC TCTAGGAAAC AGTTTAAGAA ATCATTGGCC 2251 CCTTCCCAGC ACATTGAATG GGTAAGCAGA CAGGCCATGA TTTAGTTGGC ACOCOTO CONTROL ACTIONAL ACTIONAL ACACOTO CALLACTER ACACOTO CALLAC 20 CTGCTTTAGG ATGACACAAT GAATAACACC TAGTCATATAGA AATCAGTCTC 2401 TCTGGTTTGT TTTGTATTAT GTTGTACATC ATTAAAGATC TAAATACAAA 2451 GGATATACAG TCTTGAATCT AAAATAATTT GCTAACTATT TTGATTCTTC 25DL AGAGAGAACT ACTAATAAAA ATCTAAAAGG TAAAAAAAA AAAAAAAAA 2551 A 25 **BLAST Results** ------30 No BLAST result Medline entries 35 No Medline entry 40 Peptide information for frame 1 ORF from 37 bp to 1542 bp; peptide length: 502 Category: putative protein 45 Classification: unclassified 1 MASKKLGADF HGTFSYLDDV PFKTGDKFKT PAKVGLPIGF SLPDCLQVVR 51 EVQYDFSLEK KTIEWAEEIK KIEEAEREAE CKIAEAEAKV NSKSGPEGDS JOJ KWZFZKTHZT ATMPPPINPI LAZLQHNZIL TPTRVZZZAT KQKVLSPPHI 50 151 KADFNLADFE CEEDPFDNLE LKTIDEKEEL RNILVGTTGP IMAQLLDNNL 201 PRGGSGSVLQ DEEVLASLER ATLDFKPLHK PNGFITLPQL GNCEKMSLSS 251 KVSLPPIPAV SNIKSLSFPK LDSDDSNQKT AKLASTFHST SCLRNGTFQN TOYZIVZIZA MWZZTIARME TEZMINIAZI DITHHENIEZ AZZWIZANIZ LOE

351 EESSPPNTGP TVTPPNFSVS QVPNMPSCPQ AYSELQMLSP SERQCVETVV 401 NMGYSYECVL RAMKKKGENI EQILDYLFAH GQLCEKGFDP LLVEEALEMH

451 QCSEEKMMEF LQLMSKFKEM GFELKDIKEV LLLHNNDQDN ALEDLMARAG

55

503 AS

BLASTP hits

5 No BLASTP hits available Alert BLASTP hits for DKFZphtes3_10nl0, frame 1 No Alert BLASTP hits found 10 Pedant information for DKFZphtes3_10nl0, frame 1 Report for DKFZphtes3_lOnlO.l 15 ELENGTHD 502 EMWI 55083-78 [pI] 5.02 20 EBLOCKSI PROLOBED EBFOCKZI BF0730PB EKWI All_Alpha EKWI LOW_COMPLEXITY 8-57 % 25 MASKKLGADFHGTFSYLDDVPFKTGDKFKTPAKVGLPIGFSLPDCLQVVREVQYDFSLEK SEQ SEG PRD 30 SEQ KTIEWAEEIKKIEEAEREAECKIAEAEAKVNSKSGPEGDSKMSFSKTHSTATMPPPINPI SEG PRD LASLQHNSILTPTRVSSSATKQKVLSPPHIKADFNLADFECEEDPFDNLELKTIDEKEEL SEQ 35 SEG PRD RNILVGTTGPIMAQLLDNNLPRGGSGSVLQDEEVLASLERATLDFKPLHKPNGFITLPQL SEQ SEG 40 PRD SEQ GNCEKMSLSSKVSLPPIPAVSNIKSLSFPKLDSDDSNQKTAKLASTFHSTSCLRNGTFQN SEG PRD 45 SEQ ZLKPZTQZZAZELNGHHTLGLZALNLDZGTEMPALTZZQMPZLZVLZVCTEESZPPNTGP SEG ------xxxxx PRD SEQ 50 TVTPPNFZVZQVPNMPZCPQAYSELQMLZPZERQCVETVVNMGYZYECVLRAMKKKGENI SEG PRD CCCCCCCCCCCCcchhhhhhhhccccchhhhhhhhcccchh SEQ EQILDYLFAHGQLCEKGFDPLLVEEALEMHQCSEEKMMEFLQLMSKFKEMGFELKDIKEV 55 SEG PRD SEQ LLLHNNDQDNALEDLMARAGAS

5 (No Prosite data available for DKFZphtes3_10n10.1)

(No Pfam data available for DKFZphtes3_lOnlO.l) DKFZphtes3_llal7

10

25

group: transmembrane protein

DKFZphtes3_llal7 encodes a novel 428 amino acid protein without similarity to known proteins.

The novel protein contains 2 transmembrane regions and one leucine zipper. The protein is ubiquitously expressed with higher abundance in stomach, brain and testis.

20 No informative BLAST results; No predictive prosite; pfam or SCOP motife.

The new protein can find application in studying the expression profile of testis-specific genes and as a new marker for testicular cells.

unknown protein

30 Pedant: TRANSMEMBRANE 2 perhaps differential polyadenylation

Sequenced by Qiagen

35 Locus: unknown.

Insert length: 2591 bp
Poly A stretch at pos. 2570, polyadenylation signal at pos. 2548

L CTCTCCTGCG CCCTCTGGAG GAAGTGAGAA GAGTCAGTCC CACCCAGCTG

51 CCGCCTGGTA TCTGGGCTCC AGGCCACCGA GTATTTGGCC CCCAGCCACG

101 GAGCCCTTAG CACACACCTC CCCCACAGGT CCTGGAGATG TGGCTGAGCT

151 ACCTGCAGCC GTGGCGGTAC GCGCCTGACA AGCAGGCTCC GGGCAGCGAC

45 201 TCCCAGCCCC GGTGTGTGTC GGAGAAATAGG GCACCCTTTG TCCAGGAGAA

251 CCTGCTGATG TACACCAAGT TGTTTGTGGG CTTTCTGAAC CGCGCGTCC

301 GCACAGACCT GGTCAGCCCC AAGCACGCGC TCATGGTGTT CCGAGTGGCC

351 AAAGTCTTTG CCCAGCCCAA CCTGGCTGAG ATGATTCAGA AAGGTGAGCA

401 GCTATTCCTG GAGCCAGAGC TGGTCATCCC CCACCGCCAG CACCGACTCT

50 451 TCACGGCCCC CACATTCACT GGGAGCCTTCC TGTCACCCTG GCCACCAGCG

501 GTCACTGATG CCTCCTTCAA GGTGAAGAGC CACGTCTACA GCCTGGAGGG

551 CCAGGACTGC AAGTACACCC CGATGTTTGG GCCCGAGGCC CGCACCCTGG

601 TCCTGCGCCT CGCTCAGCTC ATCACACAGG CCAAACACAC AGCCAAGTCC

651 ATCTCCGACC AGTGTGCGGA GAGCCCGGCT GGCCACTCCT TCCTCTCATG

55 701 GCTGGGCTTT AGCTCCATGG ACACCAATGG CTCCTACACA GCCAACGACC

751 TGGACGAGAT GGGGCAAGAC AGTGTCCGGA AGACAGATGA ATACCTGGAG

801 AAGGCCCTGG AGTACCTGCC TGGGCACCCAC CCCAGGCG AAGCGCAGCT

851 CAGGCAGTTC ACACTCGCCT TGGGCACCCAC CCCAGGGA AATGGAAAAA

POL AGCAACTCCC CGACTGCATC GTGGGTGAGG ACGGACTCAT CCTTACGCCC 951 CTGGGGCGGT ACCAGATCAT CAATGGGCTG CGAAGGTTTG AAATTGAGTA BDDB CCAGGGGGAC CCGGAGCTGC AGCCCATCCG GAGCTATGAG ATCGCCAGCT 1051 TGGTCCGCAC ACTCTTTAGG CTGTCGTCTG CCATCAACCA CAGATTTGCA LIDL GGACAGATGG CGGCTCTGTG TTCCCGGGAT GACTTCCTCG GCAGCTTCTG 5 1151 TCGCTACCAC CTCACAGAAC CTGGGCTGGC CAGCAGGCAC CTGCTGAGCC 1201 CTGTGGGGCG GAGGCAGGTG GCCGGCCACA CCCGCGGCCC CAGGCTCAGC 1251 CTGCGCTTCC TGGGCAGTTA CCGGACGCTG GTCTCGCTGC TGCTGGCCTT LEGAL CTTCGTGGCC TCTCTGTTCT GCGTCGGGCC CCTCCCATGC ACGCTGCTGC 1351 TCACCCTGGG CTATGTCCTC TACGCCTCTG CCATGACACT GCTGACCGAG 10 1401 CGGGGGAAGC TGCACCAGCC CTGAAGGTGT CAGCTGCCTT CAGAGCAGGC 1451 TGGAGGGATT TGCCACACAG CCCCACCCTT GGGCTGAGAG GACCTGGGAA 1501 GCCCCTCCAG GAGGGAACAC GGTCATCCTC GGGCTTCTGG AGCGGGGTTC 1551 CTGCAGCCGC AGAGGCATCT GGAGGAAACG CAACCAAGAA AGGAAGGCAG 1601 GTGGGCCCCA GCAAAGGAGT AGCTGCCAGG GCTCAACAGC TACGCTCTGT 15 1651 GACAGCGCAG AGCTCAGCGC CGGCCTTTCC CTCCCTCCGC CAAGGACTCA 17D1 CGGCCAAGCC AGCTCTCGGG GCCTTTTTTC CAGTGCCCAT TTGGCTACTC 1751 TGCTGCACCA AGCTTGGGAG CCAGCCTGCC AACAGCCACC TGGGCCTGGC LADL CTCCCCACTG GCTGGCCTTG AGGTTGGCAG AGTGGGTTGT GGCGCTTCCT LASL CTCTCTGTGT GGGACCAGGA CAGTGGCTTA AGTCTCCACT CCAGGAAAGA 20 1901 ATCAAAGTTT CTAGAGTTGT GAGAAAACCA GAGAGTGGCT GTCCTGATTC 1951 TTCACTGTGA GGGGCGTTCT TCATGTTCTC CCAGCTGTTC CAAGACTGGG 2001 CCGTAGAATT CCATGTTTCA GGAGCCTAAG ACCCTCCCAG AGCCCAGGGG 2D51 CTTCACCGCA GACCCCAAGC CATTGAGCAC ATCACCCAAA GCAGTGGCCA 2101 ACATCGCGGA CCCCTGTGCC TTGTCACAGA TGGGTGCTGG TCCTCAGGCG 25 2151 TTGGGGACAC TGCTGGGTCG ATGGGGTCGG ATTCTGCCAG TTTCTGCTCT 2201 GCAGCCAAAG ATGGTCAGAA GCATTGTCAC TTCAGTAACA TCAAGTGCTC 2251 AAAGACATGG CAACCGTTCA GTGGTACTTA AGTATTCAAA ATATACAACT 2301 ACAGATTCTC TGACAGAAAC CAGCACGGGG TCTTCACCTT CATTCACCCC 2351 ACAGGCGACA TGCGAGGGAG AACAGCATCT CAGTGGTGAT TTCCAAACCA 30 2401 AGCCTTTGTT TTCGGTGTGG GGTTTTGGGG GTTTGCTTTA ATGTTTTTGA 2451 AATTGTAAAT GTTGGGCTTT TTATTTTGAT GTAAACTGAG AATAATGGCA 2501 TTTTAGGGCC TGTGACCAAA AATGAAGCTT GTAACGACCA TGGATCTGAA 2551 TAAACATGTC CTTGCTTCTG AAAAAAAAA AAAAAAAAA A 35

BLAST Results

40 Entry AFD52134 from database EMBLNEW:
Homo sapiens clone 23585 mRNA sequence.
Score = 5765, P = 2.9e-254, identities = 1155/1156
3' UTR

45

Medline entries

50 No Medline entry

Peptide information for frame 3

55

ORF from 138 bp to 1421 bp; peptide length: 428 Category: putative protein

Classification: Transmembrane proteins unclassified Prosite motifs: LEUCINE_ZIPPER (404-425)

1 MWLSYLQPWR YAPDKQAPGS DSQPRCVSEK WAPFVQENLL MYTKLFVGFL
51 NRALRTDLVS PKHALMVFRV AKVFAQPNLA EMIQKGEQLF LEPELVIPHR
101 QHRLFTAPTF TGSFLSPWPP AVTDASFKVK SHVYSLEGQD CKYTPMFGPE
151 ARTLVLRLAQ LITQAKHTAK SISDQCAESP AGHSFLSWLG FSSMDTNGSY
201 TANDLDEMGQ DSVRKTDEYL EKALEYLRQI FRLSEAQLRQ FTLALGTTQD
251 ENGKKQLPDC IVGEDGLILT PLGRYQIING LRRFEIEYQG DPELQPIRSY
301 EIASLVRTLF RLSSAINHRF AGQMAALCSR DDFLGSFCRY HLTEPGLASR
351 HLLSPVGRRQ VAGHTRGPRL SLRFLGSYRT LVSLLLAFFV ASLFCVGPLP
401 CTLLLTLGYV LYASAMTLLT ERGKLHQP

15

BLASTP hits

No. BLASTP hits available

20

35

Alert BLASTP hits for DKFZphtes3_llal?, frame 3

No Alert BLASTP hits found

25 Pedant informa

Pedant information for DKFZphtes3_llal7, frame 3

Report for DKFZphtes3_llal7.3

30

ELENGTHD 428

EMWD 48274.93

Epid 8.92

EPROSITED LEUCINE

EPROSITED LEUCINE_ZIPPER L

EKWD TRANSMEMBRANE 2

EKUD LOW_COMPLEXITY 7.48 %

MWLSYLQPWRYAPDKQAPGSDSQPRCVSEKWAPFVQENLLMYTKLFVGFLNRALRTDLVS ZEQ 40 SEG PRD MEM PKHALMVFRVAKVFAQPNLAEMIQKGEQLFLEPELVIPHRQHRLFTAPTFTGSFLSPWPP SEQ 45 SEG PRD MEM **AVTDASFKVKSHVYSLEGQDCKYTPMFGPEARTLVLRLAQLITQAKHTAKSISDQCAESP** SEQ 50 ZEG PRD MEM AGHSFLSWLGFSSMDTNGSYTANDLDEMGQDSVRKTDEYLEKALEYLRQIFRLSEAQLRQ SEQ 55 SEG PRD MEM

	W	O 01/98454	PCT/IB01/02050
	SEQ	FTLALGTTQDENGKKQLPDCIVGEDGLILTPLGRYQIINGLRRFE	
_	PRD MEM	hhhhhhccccccccccceeecccccccccceeeecchhhhh	
5	SEQ SEG	EIASLVRTLFRLSSAINHRFAG@MAALCSRDDFLGSFCRYHLTEP	
10	PRD Mem	hhhhhhhhhhhhhhhhhhhhhhhhhhhccccceeeeeeccc	
10	SEQ	VAGHTRGPRLSLRFLGSYRTLVSLLLAFFVASLFCVGPLPCTLLL*	
15	PRD MEM	cccccccccccchhhhhhhhhhhhhhcccccchhhhhh	
10	SEQ	ERGKLHQP	
20	PRD MEM	hhhcccc	
			_
		Prosite for DKFZphtes3_11a17.	3
25	002q	DO29 404->426 LEUCINE_ZIPPER F	P\$000029
	(No	Pfam data available for DKFZphtes3_11a17.3)	·

DKFZphtes3_11c22

5 group: signal transduction

DKFZphtes3_llc22 encodes a novel 482 amino acid protein with partial similarity to mouse PC326.

The novel protein contains WD-repeats. WD-repeat proteins are known as regulatory elements in a large variety of pathways. The repeats form a propeller like strcture, which serves as a platform for protein/protein interaction. The new protein is ubiquitously expressed, indicating that it takes an essential regulatory function in the cell.

The new protein can find application in modulating/blocking of regulatory pathways.

20

25

similarity to mouse PC326

perhaps complete cds.
contains WD-Repeats: cf. BLASTX-S37694
perhaps differential polyadenylation

Sequenced by Qiagen

Locus: /map="lq23.2-24.3"

30

Insert length: 1952 bp

Poly A stretch at pos. 1932, polyadenylation signal at pos. 1912

35 L GAAGCAAGTG AGGTTGCACA AAGCAATAGA GGACGAGGAA GATCTCGACC 51 CAGAGGTGGA ACAAGTCAAT CAGATATTTC AACTCTTCCT ACGGTCCCAT LOL CAAGTCCTGA TTTGGAAGTG AGTGAAACTG CAATGGAAGT AGATACTCCA 151 GCTGAACAAT TTCTTCAGCC TTCTACATCC TCTACAATGT CAGCTCAGGC **ZDL TCATTCGACA TCATCTCCCA CAGAAAGCCC TCATTCTACT CCTTTGCTAT** 25% CTTCTCCAGA TAGTGAACAA AGGCAGTCTG TTGAGGCATC TGGACACCAC 40 **BOL ACACATCATC AGTCTGATTC TCCTTCTTCT GTGGTTAACA AACAGCTCGG** 351 ATCCATGTCA CTTGACGAGC AACAGGATAA CAATAATGAA AAGCTGAGCC 4Db CCAAACCAGG GACAGGTGAA CCAGTTTTAA GTTTGCACTA CAGCACAGAA 451 GGAACAACTA CAAGCACAAT AAAACTGAAC TTTACAGATG AATGGAGCAG 45 501 TATAGCATCA AGTTCTAGAG GAATTGGGAG CCATTGCAAA TCTGAGGGTC 551 AGGAGGAATC TTTCGTCCCA CAGAGCTCAG TGCAACCACC AGAAGGAGAC LOL AGTGAAACAA AAGCTCCTGA AGAATCATCA GAGGATGTGA CAAAATATCA **L51 GGAAGGAGTA TCTGCAGAAA ACCCAGTTGA GAACCATATC AATATAACAC** 7DL AATCAGATAA GTTCACAGCC AAGCCATTGG ATTCCAACTC AGGAGAAAGA 751 AATGACCTCA ATCTTGATCG CTCTTGTGGG GTTCCAGAAG AATCTGCTTC 50 BD1 ATCTGAAAAA GCCAAGGAAC CAGAAACTTC AGATCAGACT AGCACTGAGA BSL GTGCTACCAA TGAAAATAAC ACCAATCCTG AGCCTCAGTT CCAAACAGAA 901 GCCACTGGGC CTTCAGCTCA TGAAGAAACA TCCACCAGGG ACTCTGCTCT 951 TCAGGACACA GATGACAGTG ATGATGACCC AGTCCTGATC CCAGGTGCAA 55 LDDL GGTATCGAGC AGGACCTGGT GATAGACGCT CTGCTGTTGC CCGTATTCAG 1051 GAGTTCTTCA GACGGAGAAA AGAAAGGAAA GAAATGGAAG AATTGGATAC ILD TTTGAACATT AGAAGGCCGC TAGTAAAAAT GGTTTATAAA GGCCATCGCA 1151 ACTCCAGGAC AATGATAAAA GAAGCCAATT TCTGGGGTGC TAACTTTGTA

WO 01/98454



5	1201 1251 1301 1351 1401	TGAGCATTTG AGCCACATCC ATAAAGATCT	CTGACTGTGG ATGCTTCTGG GTTTGACCCA GGTCACCATT GTTATAACTC	AAGCTGATAA ATTTTAGCCT AGAAGAGTCA	TCATGTGGTA CATCTGGCAT AGGATTTTTA	GGCACACTGC AACTGCCTGC AGATTATGAC ACCGAAAACT GAAACTAGAA
5	1451	ACACCATTAC	AGTTCCAGCC GAGCTGACCG	TCTTTCATGT	TGAGGATGTT	GGCTTCACTT
	1551	TCAAGAGAAT		ATGAGGAATA	ATAAACTCTT	TTTGGCAAGC
10	1651 1701		AGTGCAATTT	TAAGGTTATG		TTTTTCCCTT
	1751 1801	GGAGATTGTA TGTATGAGGA	TAAAACAAAA		GTTTTTAAAA	
15	1851		GCTTGGATCA CAAATAAATT			
12	1951	AA	CAAATAAATT	ICIACACIIG	ССААААААА	мананана

BLAST Results

20

Entry HS702Jl9 from database EMBL:
Human DNA sequence *** SEQUENCING IN PROGRESS *** from clone
702Jl9

25 Score = 2043, P = 5.8e-252, identities = 425/445 10 exons matching Bp 316-1932

Entry HS536148 from database EMBL: human STS WI-6347.

30 Score = 1203, P = 1.5e-47, identities = 247/252

Entry HS703H14 from database EMBLNEW:
Human DNA sequence from clone 703H14 on chromosome 1q23.2-24.3
Score = 1307, P = 1.1e-51, identities = 263/265
2 exons matching Bp 1-316

Medline entries

40

.35

93026383:

Bergsagel PL, Timblin CR, Eckhardt L, Laskov R, Kuehl WM.; Sequence and

expression of a murine cDNA encoding PC326, a novel gene expressed in plasmacytomas but not normal plasma cells.

Oncogene
1992 Oct;7(10):2059-64

50

Peptide information for frame 1

55

ORF from 133 bp to 1578 bp; peptide length: 482 Category: similarity to known protein Classification: Protein management

Prosite motifs: MYB_1 (410-418)

1 MEVDTPAEQF LQPSTSSTMS AQAHSTSSPT ESPHSTPLLS SPDSEQRQSV
5 S1 EASGHHTHAQ SDSPSSVVNK QLGSMSLDEQ QDNNNEKLSP KPGTGEPVLS
101 LHYSTEGTTT STIKLNFTDE WSSIASSSRG IGSHCKSEGQ EESFVPQSSV
151 QPPEGDSETK APEESSEDVT KYQEGVSAEN PVENHINITQ SDKFTAKPLD
201 SNSGERNDLN LDRSCGVPEE SASSEKAKEP ETSDQTSTES ATNENNTNPE
251 PQFQTEATGP SAHEETSTRD SALQDTDDSD DDPVLIPGAR YRAGPGDRRS
10 301 AVARIQEFFR RRKERKEMEE LDTLNIRRPL VKMVYKGHRN SRTMIKEANF
351 WGANFVMSGS DCGHIFIWDR HTAEHLMLLE ADNHVVNCLQ PHPFDPILAS
401 SGIDYDIKIW SPLEESRIFN RKLADEVITR NELMLEETRN TITVPASFML

15

BLASTP hits

No BLASTP hits available

20

Alert BLASTP hits for DKFZphtes3_11c22, frame 1

TREMBLNEW: HSObb31_1 gene: "H326"; Human (H326) mRNA, complete cds., N

25 = 1, Score = 278, P = 4e-22

PIR:S37b94 gene PC32b protein - mouse, N = 1, Score = 2b5, $P = 2 \cdot 9e-20$

- 30 PIR:TO5676 hypothetical protein F20Ml3.40 Arabidopsis thaliana, N = l, Score = 240, P = 6.3e-l8
- 35 >TREMBLNEW:HSObb31_1 gene: "H326"; Human (H326) mRNA, complete cds-Length = 597

HSPs:

40

Score = 278 (41.7 bits), Expect = 4.0e-22, P = 4.0e-22 Identities = 63/148 (42%), Positives = 94/148 (63%)

Querv: 335 YKGHRNSRTMIKEANFWG--

45 ANFVMSGSDCGHIFIWDRHTAEHLMLLEADNH-VVNCLQP 391

YKGHRN+ T +K NF+G + FV+SGSDCGHIF+W++ + + + + E D

VVNCL+P

Sbict: 428 YKGHRNNAT-

VKGVNFYGPKSEFVVSGSDCGHIFLWEKSSCQIIQFMEGDKGGVVNCLEP 486

50

Query: 392 HPFDPILASSGIDYDIKIWSPLEESRIFNRKLADEVITRNELMLEE-TRNTITVPASFML 450

HP P+LA+SG+D+D+KIW+P E+ L D VI +N+ +E + +

+ S ML

55 Sbjct: 487 HPHLPVLATSGLDHDVKIWAPTAEASTELTGLKD-VIKKNKRERDEDSLHQTDLFDSHML 545

Query: 451 RMLASLNHIRADRLEGD-RSEGSGGENENEDE 481

L ++H+R R R G G + + DE Sbjct: 546 WFL--MHHLRQRRHHRRWREPGVGATDADSDE 575

5 Pedant information for DKFZphtes3_llc22, frame l

Report for DKFZphtes3_11c22.1

			•
10		•	
	FIFN	CHTO	482
	EMMI		53470 - 92
			4.72
	[pI]		• • •
	EHOM		PIR:TO4961 hypothetical protein T12J5.10 -
15			is thaliana 2e-22
	EFUN	CATI	30.09 organization of intracellular transport vesicles
		ES.	cerevisiae, YDL145cI 4e-05
	EFUN		OB-O7 vesicular transport (golgi network, etc.) [5.
			en YDL145c1 4e-05
20			
20	EFUN		99 unclassified proteins
	2e-0		
	EZUP	FAMI	WD repeat homology 4e-21
	EPRO	EJTIZ	MYB_1 1
	EKWI		Alpha_Beta
25	EKWI		LOW_COMPLEXITY 17.01 %
23	E 17 W 3		Eva_com cextri
		M=112	70.15.45.1 ADDROGRADA A LUCEDO DE CONTROL DE
	SEQ	MEAN	PAEQFLQPSTSSTMSAQAHSTSSPTESPHSTPLLSSPDSEQRQSVEASGHHTHHQ
	SEG	• • • • •	······
30	PRD	cccc	ceeeeeccccceeeeeecccccccccceeeccccchhhhhh
	SEQ	SDSPS	SZYVNKQLGSMSLDEQQDNNNEKLSPKPGTGEPVLSLHYSTEGTTTSTIKLNFTDE
	SEG		, and the same of
	PRD	6666	
35	ГКУ		
33	054		A PORT COLLEGE CAPPORTED ADDITIONS OF THE COLUMN ASSESSMENT AND ADDITIONS OF THE COLUMN ASSESSMENT
	SEQ		ASSSRGIGSHCKSEG@EESFVP@SSV@PPEGDSETKAPEESSEDVTKY@EGVSAEN
	SEG		(XXXXXXX
	PRD	cccc	
			·
40	SEQ	PVENH	INITQSDKFTAKPLDSNSGERNDLNLDRSCGVPEESASSEKAKEPETSDQTSTES
	SEG		······××××××××××××××××××××××××××××××××
	PRD	55500	eeeeecccccccccccccccccchhhhhhhhcccccccc
	LI/A	cccee	seasecccccccccccccccccccccccuuuuuuucccccccc
	55.4		NEW DEC AT A TEXT OF THE AT TH
	SEQ	ATNEN	INTNPEPQFQTEATGPSAHEETSTRDSALQDTDDSDDDPVLIPGARYRAGPGDRRS
45	SEG		(XXX
	PRD	cccc	CCCCCCEEEECCCCCCCCCCCCCCCCCCCCCCCCCCCCC
	SEQ	AVART	QEFFRRRKERKEMEELDTLNIRRPLVKMVYKGHRNSRTMIKEANFWGANFVMSGS
	SEG		•XXXXXXXXXXXXX
50	PRD		hhhhhhhhhhhhhhhhhhhccccceeeeecccccceeeeccc
50	עאים	nnnnu	MINIMINION NATIONAL DE LA CONTRACTOR DE
	SEQ	DCGHI	FIWDRHTAEHLMLLEADNHVVNCLQPHPFDPILASSGIDYDIKIWSPLEESRIFN
	SEG		
	PRD	cccee	eeeecchhhhhhhhcccceeeecccccceeeccccceeeeccchhhhhh
55			
	SEQ	BKI YI	EVITRNELMLEETRNTITVPASFMLRMLASLNHIRADRLEGDRSEGSGGENENED
		IVERA	EATIMATEMETE IMMITIALWOINTMOTHATIMATIVE CONVOEROR MENENED
	SEG		
	PRD	nncnn	hhhhhhhhhhhhhcceeecchhhhhhhhhhhhhhhhhhcccccc
			200



SEQ EE SEG .. PRD cc

5

Prosite for DKFZphtes3_11c22.1

7E00029 01

· 410->419 MYB_1

PD0C00037

(No Pfam data available for DKFZphtes3_llc22.l)

WO 01/98454

DKFZphtes3_11d21



5 group: signal transduction

DKFZphtes3_lld2l encodes a novel 922 acid protein and contains the full coding sequence of the human Nedd-4-like ubiquitinprotein ligase.

The novel protein contains four WW domains. The WW/rsp5/WWP domain has been shown to bind proteins with particular prolinemotifs, and thus resembles somewhat SH3 domains. It is frequently associated with other domains typical for proteins in signal transduction processes. There is also a ubiquitin-protein ligase activity reported. The protein is believed to play an important role in protein-degradation pathways.

The new protein can find application in diagnosis of diseases due to unnormal protein degradation like muscular dystrophy or multiple sclerosis as well as in modulating the half life of specific proteins and in expression profiling.

25 similarity to Nedd-4-like ubiquitin-protein ligase (Homo sapiens)
Sequenced by Qiagen

Locus: unknown

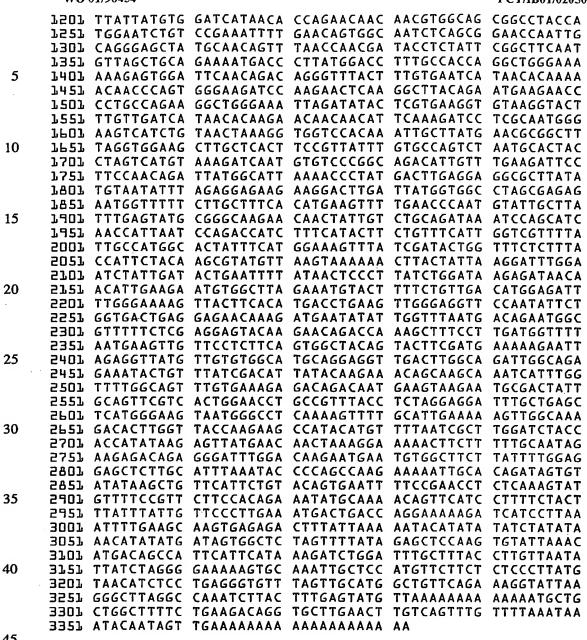
30

Insert length: 3382 bp

Poly A stretch at pos. 3362, polyadenylation signal at pos. 3345

35 1 ATTTTGGGAC ATGGCCACTG CTTCACCAAG GTCTGATACT AGTAATAACC 51 ACAGTGGAAG GTTGCAGTTA CAGGTAACTG TTTCTAGTGC CAAACTTAAA 101 AGAAAAAGA ACTGGTTCGG AACAGCAATA TATACAGAAG TAGTTGTAGA 151 TGGAGAAATT ACGAAAACAG CAAAATCCAG TAGTTCTTCT AATCCAAAAT 201 GGGATGAACA GCTAACTGTA AATGTTACGC CACAGACTAC ATTGGAATTT 40 251 CAAGTTTGGA GCCATCGCAC TTTAAAAGCA GATGCTTTAT TAGGAAAAGC BOL AACGATAGAT TTGAAACAAG CTCTGTTGAT ACACAATAGA AAATTGGAAA 351 GAGTGAAAGA ACAATTAAAA CTTTCCTTGG AAAACAAGAA TGGCATAGCA 401 CAAACTGGTG AATTGACAGT TGTGCTTGAT GGATTGGTGA TTGAGCAAGA 451 AAATATAACA AACTGCAGCT CATCTCCAAC CATAGAAATA CAGGAAAATG 45 501 GTGATGCCTT ACATGAAAAT GGAGAGCCTT CAGCAAGGAC AACTGCCAGG 551 TTGGCTGTTG AAGGCACGAA TGGAATAGAT AATCATGTAC CTACAAGCAC LOS TCTAGTCCAA AACTCATGCT GCTCGTATGT AGTTAATGGA GACAACACAC L51 CTTCATCTCC GTCTCAGGTT GCTGCCAGAC CCAAAAATAC ACCAGCTCCA 701 AAACCACTCG CATCTGAGCC TGCCGATGAC ACTGTTAATG GAGAATCATC 50 751 CTCATTTGCA CCAACTGATA ATGCGTCTGT CACGGGTACT CCAGTAGTGT BOL CTGAAGAAAA TGCCTTGTCT CCAAATTGCA CTAGTACTAC TGTTGAAGAT 851 CCTCCAGTTC AAGAAATACT GACTTCCTCA GAAAACAATG AATGTATTCC 901 TTCTACCAGT GCAGAATTGG AATCTGAAGC TAGAAGTATA TTAGAGCCTG
951 ACACCTCTAA TTCTAGAAGT AGTTCTGCTT TTGAAGCAGC CAAATCAAGA
1001 CAGCCAGATG GGTGTATGGA TCCTGTACGG CAGCAGTCTG GGAATGCCAA 55 1051 CACAGAAACC TTGCCATCAG GGTGGGAACA AAGAAAAGAT CCTCATGGTA
1101 GAACCTATTA TGTGGATCAT AATACTCGAA CTACCACATG GGAGAGACCA 1151 CAACCTTTAC CTCCAGGTTG GGAAAGAAGA GTTGATGATC GTAGAAGAGT





BLAST Results

50 No BLAST result

Medline entries

55

45

Pirozzi Ga McConnell SJa Uveges AJa Carter JMa Sparks ABa Kay BKa Fowlkes DM.: Identification of novel human WW domain-containing



WO 01/98454



proteins
by cloning of ligand targets J Biol Chem 1997 Jun
6:272(23):14611-6

5

Peptide information for frame 2

10

ORF from 11 bp to 2776 bp; peptide length: 922

Category: known protein

Classification: Protein management Prosite motifs: WW_DOMAIN_1 (355-380)

15 WW_DOMAIN_1 (387-412) WW_DOMAIN_1 (462-487) WW_DOMAIN_1 (502-527)

1 MATASPRSDT SNNHSGRLQL QVTVSSAKLK RKKNWFGTAI YTEVVVDGEI 51 TKTAKSSSS NPKWDEQLTV NVTPQTTLEF QVWSHRTLKA DALLGKATID 101 LKQALLIHNR KLERVKEQLK LSLENKNGIA QTGELTVVLD GLVIEQENIT 20 151 NCSSSPTIEI GENGDALHEN GEPSARTTAR LAVEGTNGID NHVPTSTLVQ 201 NSCCSYVUG DUTPSSPSQV AARPKUTAAR KPLASERAD TVNGESSSFA 251 PTDNASVTGT PVSEENALS PNCTSTTVED PPVQEILTSS ENNECIPSTS 301 AELESEARSI LEPDTSNSRS SSAFEAAKSR QPDGCMDPVR QQSGNANTET 25 351 LPSGWERKD PHGRTYYVDH NTRTTTWERP RPLPPGWERR VDDRRRVYYV 401 DHNTRTTTWR RPTMESVRNF ERWRSRRWRL RGAMRRFNRR YLYSASMLAA 451 ENDPYGPLPP GWEKRVDSTD RVYFVNHNTK TTQWEDPRTQ GLQNEEPLPE 501 GWEIRYTREG VRYFVDHNTR TTTFKDPRNG KSSVTKGGPQ IAYERGFRWK 551 LAHFRYLCQS NALPSHVKIN VSRQTLFEDS FQQIMALKPY DLRRRLYVIF 30 LOD RGEEGLDYGG LAREWFFLLS HEVLNPMYCL FEYAGKNNYC LQINPASTIN L51 PDHLSYFCFI GRFIAMALFH GKFIDTGFSL PFYKRMLSKK LTIKDLESID 701 TEFYNSLIWI RDNNIEECGL EMYFSVDMEI LGKVTSHDLK LGGSNILVTE 751 ENKDEYIGLM TEWRFSRGVQ EQTKAFLDGF NEVVPLQWLQ YFDEKELEVM 35 ADD LCGMQEVDLA DWQRNTVYRH YTRNSKQIIW FWQFVKETDN EVRMRLLQFV 851 TGTCRLPLGG FAELMGSNGP @KFCIEKVGK DTWLPRSHTC FNRLDLPPYK 901 SYEQLKEKLL FAIEETEGFG QE

40

BLASTP hits

No BLASTP hits available

45

Alert BLASTP hits for DKFZphtes3_11d21, frame 2

No Alert BLASTP hits found

50

Pedant information for DKFZphtes3_11d21, frame 2

Report for DKFZphtes3_11d21.2

55

ELENGTHD 925

EMW3 105650.58

EpI3 5-60

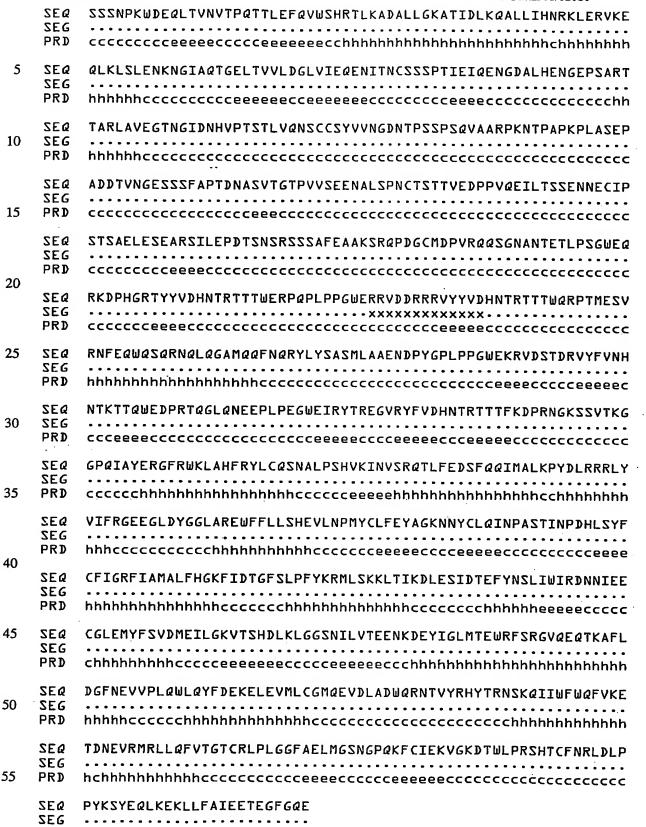


WO 01/98454



TREMBL:HSU96113_1 product: "WWP1"; Homo sapiens Nedd-4-like ubiquitin-protein ligase WWPL mRNA, partial cds. 0.0 EFUNCATI 30.02 organization of plasma membrane ES. cerevisiae, YER125wl le-149 [FUNCAT] 11.01 stress response ES. cerevisiae, YER125wl le-149 **EFUNCAT** 06.13.01 cytoplasmic degradation ES. cerevisiae, YER125wl le-149 EFUNCATI 03.10 sporulation and germination ES. cerevisiae, 10 YER125wl le-149 EFUNCATI 06.07 protein modification (glycolsylation, acylation, myristylation, palmitylation, farnesylation and processing) ES. cerevisiae, YERL25wl le-149 **EFUNCATI** 03.22 cell cycle control and mitosis ES. cerevisiae, 15 YDR457wl le-78 EFUNCATE 99 unclassified proteins ES. cerevisiae, YJRD3bcJ 7e-39 **EFUNCATI** 30.03 organization of cytoplasm ES. cerevisiae, YKLOlocJ 8e-21 20 EFUNCATI 30.10 nuclear organization ES. cerevisiae, YKLOl2wl ье-05 **EFUNCATI** 04.05.03 mrna processing (splicing) ES. cerevisiae, YKLOl2wl be-05 **EFUNCATI** 30.01 organization of cell wall ES. cerevisiae. 25 YIROL9cl 3e-04 [FUNCAT] 30.90 extracellular/secretion proteins ES. cerevisiae, YIRO19cl 3e-04 EFUNCATI 01.05.01 carbohydrate utilization ES. cerevisiae, YIRO19c1 3e-04 30 EBLOCKSI BPD3746E EBLOCKZJ BPO37616 EBF0CKZ] BLOO514E Fibrinogen beta and gamma chains C-terminal domain proteins [BLOCKZ] PR007318 35 [BLOCK2] BP01566C **EBFOCKZ** BLD1159 WW/rsp5/WWP domain proteins [Brock2] PR00403B **EBFOCKZ** PR00403A **EBFOCKZ** PF00632B 40 EBLOCK23 PF00632A **EECI** 6.3.2.19 Ubiquitin--protein ligase le-151 **EPIRKU** ligase le-151 [PIRKW] transmembrane protein 2e-37 **EPIRKWI** leucine zipper 2e-28 45 WW repeat homology le-151 ESUPFAMI WD repeat homology 2e-28 **ESUPFAMD** ubiquitin ligase homolog le-151 **ESUPFAMI** EPROSITE WW_DOMAIN_1 [PFAM] WW/rsp5/WWP domain containing proteins 50 **EPFAMJ** C2 domain EKMI Alpha_Beta EKWI LOW_COMPLEXITY 1.41 %





PRD ccchhhhhhhhhhhhhhhhhhccccc

```
5
                          Prosite for DKFZphtes3_11d21.2
    P201159
                   358->384
                              MM_DOWAIN_J
                                                       PD0C50020
    P201159
                              WW_DOMAIN_L
                   390->416
                                                       PD0C50020
    P201159
                   465->491
                              MM_DOWAIN_T
                                                       PD0C50020
10
    PZ01159
                   505->531
                              WW DOMAIN 1
                                                       PD0C50020
                          Pfam for DKFZphtes3_11d21.2
15
    HMM NAME C2 domain
20
    *LtVrIIeARNLWkMDMnGfSDPYVKVdMdPdpkDtkKWKTkTiWNN.GL
                        L V++ +A+ +K+++G+
                                              Y +V +D++
                                                                TKT
    +++ +
    Query
                    23 ·L·QVTVSSAKLKRKKNWFGTA-IYTEVVVDGE-----
    ITKTAKZZZZZ
                   L 3
25
    MMH
                       NPVWNEEefvFedIPyPdlqrkMLRFaVWDWDRFSRBDFIGHCi*
                       NP W+ E+++ + + +
                                            L+F+VW + ++ +++6 ++
                    64 NPKWD-EQLTVN---VTPQTT--LEFQVWSHRTLKADALLGKAT
    Query
    דטד
30
    HMM_NAME WW/rsp5/WWP domain containing proteins
35
    HMM
                       *LPsGWEeHWDpsGRpWYYWNHETkTTQWEpP*
                        LPSGUE+++DP GR+ YY++H+T+TT+WE+P
                   354 LPSGWEQRKDPHGRT-YYVDHNTRTTTWERP
    Query
                                                            383
              JAB
                    415
                                 31 dkfzphtes3_11d21.2 similarity to
40
    Nedd-4-like ubiquitin-protein ligase (Homo sapiens)
      Alignment to HMM consensus:
                       *LPsGWEeHWDpsGRpWYYWNHETkTTQWEpP*
                        LP+GUE++ D+ R YY++H+T+TT+W++P
      dkfzphtes3
                   386 LPPGWERRVDDRRRV-YYVDHNTRTTTWQRP
                                                            435
45
                    490
                                 31 dkfzphtes3_11d21.2 similarity to
    Query
    Nedd-4-like ubiquitin-protein ligase (Homo sapiens)
      Alignment to HMM consensus:
    HMM
                       *LPsGWEeHWDpsGRpWYYWNHETkTTQWEpP*
50
                        LP+GWE++ D + R Y++NH+TKTTQWE+P
   Query
                   461 LPPGWEKRVDSTDRV-YFVNHNTKTTQWEDP
                                                            490
              501
    34.62
                    530
                            ŀ
                                 31 dkfzphtes3_11d21.2 similarity to
    Nedd-4-like ubiquitin-protein ligase (Homo sapiens)
55
      Alignment to HMM consensus:
                       *LPsGWEeHWDpsGRpWYYWNHETkTTQWEpP*
    Query
                        LP GWE +++ +G + Y+++H+T+TT+ ++P
                   501 LPEGWEIRYTREGVR-YFVDHNTRTTTFKDP
      dkfzphtes3
                                                            530
```

PAGE INTENTIONALLY LEFT BLANK

5 group: testis derived

DKFZphtes3_1le17 encodes a novel 573 amino acid protein without similarity to known proteins.

No informative BLAST results; No predictive prosite, pfam or SCOP motife.

The new protein can find application in studying the expression profile of testis-specific genes.

15

unknown protein

Sequenced by @iagen

20

Locus: unknown

Insert length: 2102 bp

Poly A stretch at pos. 2080, polyadenylation signal at pos. 2059

25

```
L GGCCTGGGGG GCTTCCCTGG GGGGCTTGTC GCCGGGGCCG CCTGGGCTTT

5L CAGGTCTTCC GAGGCTGACA TTCACGTTTC ATTCTGCCAC ACTCGGGAAC

LOL GGTGATCGGG GAAGCATGGG GATCCGGGAG AAGCACCCAC AAAACTAGCA

LSL TCCTCCTGGA GGAGCTCGGG AATAGGATGA GTGATAATCC ACCCAGAATG

20L GAAGTGTGTC CTTACTGTAA GAAGCCATTT AAACGATTAA AATCCCACTT

25L GCCATACTGT AAGATGATAG GATCAACCAT ACCTACTGAT CAAAAAGTTT

30L ATCAGTCCAA GCCAGCTACA CTCCCACGTG CTAAAAAGAT GAAAGGACCA

35L ATCAAAGATT TAATTAAAGC TAAAGGGAAA GAGTTAGAGA CAGAGAATGA

40L AGAAAGAAAT TCTAAGTTGG TGGTGGACAA ACCAGAACAG ACAGTGAAGA

45L CCTTTCCACT GCCAGCTGTT GGTTTGGAAA GAGCAGCTAC TACAAAGGCA

50L GATAAAGACA TCAAGAATCC AATCCAACCA TCCTTCAAAA TGTTAAAAAA

55L TACTAAACCA ATGACTACTT TCCAAGAAGA AACCAAGGCT CAGTTTTACG

60L CATCAGAGAA AACCTCTCCT AAAAGAGAAC TTGCCAAAGA TTTGCCTAAA
30
35
              LOT CATCAGAGAA AACCTCTCCT AAAAGAGAAC TTGCCAAAGA TTTGCCTAAA LST TCAGGAGAAA GTCGATGTAA TCCTTCAGAA GCTGGAGCGT CTTTACTGGT
40
              701 TGGCTCAATA GAACCTTCTT TGTCAAATCA AGATAGAAAA TATTCCTCAA
              751 CTCTACCTAA TGATGTACAA ACTACCTCTG GTGATCTCAA ATTGGACAAA
              BOL ATTGATCCCC AAAGACAGGA ACTTCTAGTA AAATTACTAG ATGTGCCTAC
           45
50
            1351 AGAATCATAA TTGTGTCCCT GATGTAAAGG CATTAATGGA GAGTCCCGAG
1401 GGACAGTTAT CTCTGGAGCC CAAATCTGAT AGTCAGTTCC AAGCATCACA
1451 CACTGGGTGC CAGAGCCCTT TATGTTCAGC CCAGCGTCAC ACTCCTCAGA
55
           1501 GCCCCTTCAC CAATCATGCT GCAGCTGCTG GCAGGAAGAC TCTTCGCAGC
           1551 TGCATGGGGC TGGAGTGGTT TCCAGAGCTC TATCCTGGTT ACCTTGGACT
```

WO 01/98454 PCT/IB01/02050 LLDL AGGGGTGTTG CCAGGGAAGC CTCAGTGTTG GAATGCAATG ACCCAGAAGC LLSL CACAACTTAT CAGTCCCCAG GGGGAAAGAC TCTCACAAGG CTGGATCAGG L7DL TGCAACACCA CCATAAGGAA GAGTGGATTC GGTGGCATCA CTATGCTCTT 1751 CACAGGATAC TTCGTCCTGT GTTGTAGCTG GAGTTTCAGA CGTCTGAAAA 1801 AATTGTGCCG ACCCCTGCCC TGGAAGAGCA CAGTACCTCC ATGCATTGGT 5 1851 GTGGCGAAGA CGACTGGGGA TTGCCGCTCT AAAACATGTT TGGATTAGGA 1901 AGCACGTTTA AGTAGGAGAA GCCTTCGTGA CTTCTCTCTA GTGCCTTCGT 1951 GCCCTGTGTT GCCCACTGAA TTGCCCTGTA ACACCTAAGT GTAGTGGTAG 2001 CATTAAGGGA TAGCTTTTCA GCCCTCAAGG TTATCAGGAG CATTTGTATC 2051 ACTGCTATAA ATAAAGTAGT ATCACTTGTC ATAAAAAAA AAAAAAAAA 10 5707 VV **BLAST Results** 15 -----No BLAST result 20 Medline entries No Medline entry 25 Peptide information for frame 3 30 ORF from 177 bp to 1895 bp; peptide length: 573 Category: putative protein Classification: no clue 1 MSDNPPRMEV CPYCKKPFKR LKSHLPYCKM IGSTIPTDQK VYQSKPATLP 35 51 RAKKMKGPIK DLIKAKGKEL ETENEERNSK LVVDKPEQTV KTFPLPAVGL LOL ERAATTKADK DIKNPIQPSF KMLKNTKPMT TFQEETKAQF YASEKTSPKR
LSL ELAKDLPKSG ESRCNPSEAG ASLLVGSIEP SLSNQDRKYS STLPNDVQTT
201 SGDLKLDKID PQRQELLVKL LDVPTGDCHI SPKNVSDGVK RVRTLLSNER 251 DZKGRDHLZG VPTDVTVTET PEKNTESLIL SLKMZSLGKI QVMEKQEKGL 301 TLGVETCGSK GNAEKSMSAT EKQERTVMSH GCENFNTRDS VTGKESQGER 40 45% AGRKTLRSCM GLEWFPELYP GYLGLGVLPG KPQCWNAMTQ KPQLISPQGE 501 RLSQGWIRCN TTIRKSGFGG ITMLFTGYFV LCCSWSFRRL KKLCRPLPWK 45 551 STVPPCIGVA KTTGDCRSKT CLD BLASTP hits 50 No BLASTP hits available

Alert BLASTP hits for DKFZphtes3 llel7, frame 3

55 No Alert BLASTP hits found

Pedant information for DKFZphtes3_llel7, frame 3

Report for DKFZphtes3_llel7.3

```
5
   ELENGTHD
         573
         63389.88
   EMWI
   [[q]
         9.24
         BLOODEA Zinc finger, C2H2 type, domain proteins
   EBFOCK23
   EKWI
         Alpha_Beta
10
   [KW]
         LOW_COMPLEXITY
                      7-50 %
      MSDNPPRMEVCPYCKKPFKRLKSHLPYCKMIGSTIPTDQKVYQSKPATLPRAKKMKGPIK
   SEQ
   SEG
      15
   PRD
      DLIKAKGKELETENEERNSKLVVDKPEQTVKTFPLPAVGLERAATTKADKDIKNPIQPSF
   SEQ
   SEG
   PRD
      20
      KMLKNTKPMTTFQEETKAQFYASEKTSPKRELAKDLPKSGESRCNPSEAGASLLVGSIEP
   SEQ
   SEG
   PRD
      hhhhcccccchhhhhhhhhhhcccccchhhhhhhcccc
      SLSN@DRKYSSTLPNDV@TTSGDLKLDKIDP@R@ELLVKLLDVPTGDCHISPKNVSDGVK
25
   SEQ
   SEG
   PRD
      RVRTLLSNERDSKGRDHLSGVPTDVTVTETPEKNTESLILSLKMSSLGKIQVMEKQEKGL
   SEQ
30
   SEG
      PRD
   SEQ
      TLGVETCGSKGNAEKSMSATEKQERTVMSHGCENFNTRDSVTGKESQGERPHLSLFIPRE
   SEG
      35
  PRD
      TTYQFHSVSQSSSQSLASLATTFLQEKKAEAQNHNCVPDVKALMESPEGQLSLEPKSDSQ
  SEQ
      -----×xxxxxxxxxxx------
  SEG
      PRD
40
  SEQ
      FQASHTGCQSPLCSAQRHTPQSPFTNHAAAAGRKTLRSCMGLEWFPELYPGYLGLGVLPG
  SEG
      PRD
      45
      KPQCUNAMTQKPQLISPQGERLSQGWIRCNTTIRKSGFGGITMLFTGYFVLCCSWSFRRL
  SEQ
  SEG
  PRD
      ccccccccccccchhhhhccccceeeeccccceeeecceeeecchhhhh
      KKLCRPLPWKSTVPPCIGVAKTTGDCRSKTCLD
  SEQ
50
  SEG
  PRD
     hhhcccccccccceeeeecccccccccc
  (No Prosite data available for DKFZphtes3_lle17.3)
55
  (No Pfam data available for DKFZphtes3_llel7.3)
```

5 group: testis derived

DKFZphtes3_12dl8 encodes a novel 1170 amino acid protein without similarity to known proteins.

- The EST-distribution signifies an ubiquitous expression pattern. No informative BLAST results; No predictive prosite, pfam or SCOP motife.
- The new protein can find application in studying the expression profile of testis-specific genes.

unknown protein

25

20 perhaps complete cds.

Sequenced by Qiagen

Locus: /map="136-9 cR from top of Chrl3 linkage group"

Insert length: 5469 bp

Poly A stretch at pos. 5449, polyadenylation signal at pos. 5420

30 1 AAGGACAGAG GACGAGATTT TGAACGACAA AGAGAAAAGA GAGACAAGCC 51 AAGGTCTACT TCCCCAGCAG GACAGCATCA TTCTCCTATA TCTTCTAGAC
101 ATCACTCATC TTCCTCACAA TCAGGATCAT CTATTCAAAG ACATTCTCCT 151 TCTCCTCGTC GAAAAAGAAC TCCTTCACCA TCTTATCAGC GGACACTAAC 201 TCCACCTTTA CGACGCTCTG CCTCTCCTTA TCCTTCACAT TCTTTGTCGT
251 CTCCCCAGAG AAAGCAGAGT CCTCCAAGAC ATCGCTCTCC AATGCGAGAG
301 AAAGGGAGAC ATGATCATGA ACGAACTTCA CAGTCTCATG ATCGACGCCA 35 351 CGAAAGGAGG GAAGATACTA GGGGCAAACG AGACAGAGAA AAGGACTCAA 401 GAGAAGAACG AGAATATGAA CAGGATCAGA GCTCTTCTAG AGACCACAGA 451 GATGACAGAG AACCTCGAGA TGGTCGGGAT CGGAGAGATG CCAGAGATAC 40 501 TAGGGACCGA AGGGAACTAA GAGACTCCAG AGACATGCGG GACTCAAGGG 551 AGATGAGAGA TTATAGCAGA GATACCAAAG AGAGCCGTGA TCCCAGAGAT 601 TCTCGGTCCA CTCGTGATGC CCATGACTAC AGGGACCGTG AAGGTCGAGA L51 TACTCATCGA AAGGAGGATA CATATCCAGA AGAATCCCGG AGTTATGGCC 701 GAAACCATTT GAGAGAAGAA AGTTCTCGTA CGGAAATAAG GAATGAGTCC 45 -751 AGAAATGAGT CTCGAAGTGA AATTAGAAAT GACCGAATGG GCCGAAGTAG BOL GGGGAGGGTT CCTGAGTTAC CTGAAAAGGG AAGTCGAGGC TCAAGAGGTT 851 CTCAAATTGA TAGTCACAGT AGTAATAGCA ACTATCATGA CAGCTGGGAA 901 ACTCGAAGTA GCTATCCTGA AAGAGATAGA TATCCTGAAA GAGACAACAG 951 AGATCAAGCA AGGGATTCTT CCTTTGAGAG AAGACATGGA GAGCGAGACC 50 BODD GTCGTGACAA CAGAGAGAGA GATCAAAGAC CAAGCTCACC AATTCGACAT 1051 CAGGGAAGGA ATGACGAGCT TGAGCGTGAT GAAAGAAGAG AGGAACGAAG JJDJ AGTAGACAGA GTGGATGATA GGAGAGATGA AAGGGCTAGA GAGAGAGTC 1151 GGGAACGAGA ACGAGACAGG GAGCGGGAGA GAGAGAGGGA ACGTGAACGG 1201 GATCGGGAAA GAGAAAAGA GAGAGAACTA GAAAGAGAGC GTGCTAGGGA 55 1251 ACGGGAGAGA GAAAGAGAAA AAGAGAGAG TCGTGAAAGG GATAGAGACC 1301 GAGACCACGA TCGAGAGCGG GAAAGAGAGA GGGAACGAGA CAGGGAAAA 1351 GAACGGGAAC GAGAAAGAGA AGAGAGAGA AGGGAGAGAG AGCGAGAACG 14D1 GGAGAGAGA CGAGAGCGAG AACGGGAACG AGAAAGAGCG AGAGAAAGGG

```
1451 ATAAAGAACG AGAACGCCAA AGGGATTGGG AAGACAAAGA CAAAGGACGA
      1501 GATGACCGCA GAGAAAAGCG AGAAGAGATC CGAGAAGATA GGAATCCAAG
      1551 AGATGGACAT GATGAAAGAA AATCAAAGAA GCGCTATAGA AATGAAGGGA
      ILDI GTCCCAGCCC TAGACAGTCC CCGAAGCGCC GGCGTGAACA TTCTCCGGAC
      1651 AGTGATGCCT ACAACAGTGG AGATGATAAA AATGAAAAAC ACAGACTCTT
  5
      1701 GAGCCAAGTT GTACGACCTC AAGAATCTCG TTCTCTTAGT CCCTCGCACC
      1751 TCACAGAAGA CAGACAGGGT AGATGGAAAG AGGAGGATCG TAAACCAGAA
      LBDL AGGAAAGAGA GTTCAAGGCG CTACGAAGAA CAGGAACTCA AGGAGAAAGT
      1851 TTCTTCTGTA GATAAACAGA GAGAACAGAC AGAAATCCTG GAAAGCTCAA
      ·LPOL GAATGCGTGC ACAGGACATT ATAGGACACC ACCAGTCTGA AGATCGAGAG
 10
      1951 ACATCTGATC GAGCTCATGA TGAAAACAAG AAGAAAGCAA AAATTCAAAA
      2001 GAAACCAATT AAGAAAAAGA AAGAGGATGA TGTTGGAATA GAGAGGGGTA
      2D51 ACATAGAGAC AACATCTGAA GATGGTCAAG TATTTTCACC AAAAAAAGGA
      2101 CAGAAAAGA AAAGCATTGA AAAAAACGT AAAAAATCCA AAGGTGATTC
      2151 TGATATTTCT GATGAAGAAG CAGCCCAGCA AAGTAAGAAG AAAAGAGGCC
 15
      2201 CACGGACTCC CCCTATAACA ACTAAAGAGG AATTGGTTGA AATGTGCAAT
      2251 GGTAAGAATG GTATTCTAGA GGACTCCCAG AAAAAAGAAG ATACAGCATT
      23D1 CAGTGACTGG TCTGATGAGG ATGTCCCTGA CCGTACAGAG GTGACAGAG
      2351 CAGAGCATAC TGCCACCGCC ACGACTCCTG GTAGTACCCC TTCTCCTCTA
      2401 TCTTCTCTTC TTCCTCCTCC ACCGCCTGTG GCTACTGCCA CTGCTACAAC
 20
      2451 TGTGCCTGCA ACTCTTGCTG CCACTACTGC TGCTGCCGCC ACCTCTTTCA
      25D1 GCACATCTGC CATCACTATT TCCACCTCTG CCACCCCCAC CAATACCACC
      2551 AATAATACTT TTGCCAATGA AGACTCACAC AGAAAATGCC ACAGAACACG
      2603 AGTAGAAAA GTAGAGACGC CTCACGTGAC TATAGAAGAT GCACAGCATC
      2651 GCAAGCCTAT GGATCAAAAG AGGAGCAGCA GCCTCGGGAG CAATCGGAGT
 25
      2701 AACCGTAGTC ATACGTCTGG TCGTCTTCGC TCCCCATCCA ATGATTCAGC
      2751 CCATCGAAGT GGAGATGACC AAAGTGGTCG AAAGAGAGTA CTGCACAGTG
      2801 GCTCAAGAGA TAGAGAAAAA ACAAAAAGCC TGGAAATCAC AGGAGAGAGA
      2851 AAATCTAGGA TTGATCAGTT AAAGCGTGGA GAACCCAGTC GAAGTACTTC
 30
      29D) TTCAGATCGC CAGGATTCAA GAAGCCATAG TTCAAGAAGA AGTTCTCCAG
      2951 AGTCAGATCG ACAGGTCCAT TCAAGATCTG GGTCATTTGA TAGCAGAGAC
      ADD1 AGGCTTCAAG AACGAGATCG ATATGAACAC GACAGAGAC GCGAGAGAGAGA
      3D51 GAGGAGAGAT ACGAGGCAGA GAGAATGGGA CCGAGATGCT GATAAAGATT
      BLDL GGCCACGCAA CAGGGATCGA GATAGATTGC GAGAACGAGA ACGAGAGAGA
    - 3151-GAACGAGACA AAAGGAGAGA CTTGGATAGG GAAAGAGAGA-GACTAATTTC
 35
      3201 TGATTCTGTT GAAAGGGACA GGGACAGAGA CAGAGACAGA ACTTTTGAGA
      3251 GTTCTCAAAT AGAGTCTGTG AAACGCTGTG AAGCAAAACT GGAAGGTGAA
      AATADATOO DATOTOTOA AGOTOTOA DAAADATOTA GOGAADTAO LOEE
      3351 AGAGAGATG GATAAAGATC TGGGATCTGT GCAGGGATTT GAAGATACAA
 40
      3401 ATAAATCCGA GAGAACTGAG AGTCTGGAAG CAGGAGATGA CGAGTCCAAG
      345% TTAGATGATG CACATTCATT AGGCTCTGGT GCTGGAGAAG GATACGAGCC
      3501 AATCAGTGAT GACGAACTAG ATGAAATTCT GGCAGGTGAT GCAGAAAGA
      3551 GGGAGGACCA ACAGGATGAG GAGAAGATGC CAGATCCCTT AGATGTGATA
      ADAGOACOAA BAAAAG CATCTGGTCT TATGCCAAAG CATCCAGAG AACCACGAGA
      365 GCCTGGGGCT GCACTCTTAA AATTCACACC TGGAGCTGTT ATGCTAAGAG
. 45
      3701 TTGGGATTTC TAAAAAGTTG GCAGGTTCTG AACTCTTTGC CAAAGTCAAA
      3751 GAAACATGTC AGAGACTTTT AGAAAAACCC AAAGGTAGTT TCATTTTACT
      TTAGOGOAAA TATATOTOTAA GATGCCATCT AACGATAT TATATOTATT
      3851 CTGATCTGTT CTTATGTAGC ACTTAACACT GTGTAGAAAC TATTTTTTGA
      3901 GAASTATTT TATAATCATT ATTAACCT CATGGTCAAA GTTTCTCTTT
 50
      395b AAAATTTATT TTGAGAAGAA GAGTTATCCC ACAGAAAGT TGGGAAAGA
      4DD1 GTACAATGAC CTTTTTGTAT GAAAATTACT TATTAACAGG CCAGGCGTGG
      4051 TGTTGCATGT CTGTAGTCAC AGCTACTCAG GGAGGTTGAG GCAGCAGGAT
      4101 TGCTGGAGCC CAGGAAATTG AGGCTGCAGT GAGCCATGAT TGAGCCACCA
      4151 CACTCCAACC TAGGTGACAG AGCAAGACCC TGTCTCAAAA AAAAAAAAAC 4201 AAATTAACCA ATAAGTTCTA ATATCAAAGT GCTCAGTGGT TTGCCCTTGG
 55
      AATOTOTAOT TITTATAOAT OAGACAAAA GOACOGAGA GAAGTAAATO 4254
OGTAAAACOOA TOTAOAGAT COOTTTOOO AOGTTTTAT GOOTTAAAGA LOEF
```

WO 01/98454

4351 AGGTGTGTAG GTTGAGTCTT TAACAAAGTG ATTAAGAGCT TGGTCTGTAA

4403 GGCCGGATGA TCTGGATTTC AGTAGGCACA CCACTTACTG GCTATTACTT

4451 AATCTGTGTG TTAGTGTCAT CATCTGTAAG TCAGGAATAA TCATACCACC

4501 AACTTCCTAT GGTAATTAGG AGCAAATGAG TTATTACAGG CAAAACACTT

4551 AGAACAGTTC CTGGCATATA GTAATACCCA ATAAATATTA ACTGCTACTT

4501 TGAAAATATC CTATCACGCT GATTTTTGAC CTCACTGCAG CAATTTTCAG

4551 TTATTCCAGA TTATCTAGCT TATGGATTCT GGTGGTAGGG GTTGTTTGGT

4701 TTTGGTTTTC ACTGTCTCTG TCTCATCTAG TACCTACCTT AGTTTATTTT

4751 GCAACTTACT AATACTTTAT TAATGGGGAG GGACGAGTAG ATGGTAAAAA

4751 GCAACTTACT AATACTTTAT TAATGGGGAG GGACGAGTAG ATGGTAAAAA 4801 GAAGGAAAAG GAGGTAAAAG GTGAAAGGAA CAACATTAAT TAACAATTTT 4851 ACGTCATGTC CCTGGACATA AAAGTTTAGT TAGTATTAAA TTTTTCACTA 4901 ATACAAAATA AAAAAATATT GTTTTATGAG TTTTATGAAT TCATGCCCTT

4951 CCTTTACTCT ATTAGCATAA GCAGTAAATT TTTTTATTTT AATATAGCCC
5001 AATAAACCTA GAGTATACAT GTACAAAATA CATATAATTG TTAACGTGTA
5051 TTAACCGAAA AATGACCCAA GACTTAGTTC TTGCCCTACT GTATCTGCCT

5101 TGTTTGGTTG GTTCTGTGAC CTTAAGCAAA TAACTCCTGT GAGCCTCAAT
5151 TTTATTTGTA AAGTGATGGA ATAAAACCCC TAAAATCTTA CCCACCTCTA

5201 AAGATATTG TTTCTGTGAC CTTTTGCTAG TAGCATTTCA AGTTAAAATC 5251 TGGTTTGATT TTGCTACCCA TGAAATACAG TTCGGCCCTT ACTTATTGAT 5301 GACTTAACCT AAACAGTGAA AATATGCACT GTAAAGGGTG GGGTGATGTG

5351 GCTTAACAAT CAGACTTCTT CTATTTTTGC TGCTATGGTG GTTGTATTAG
5401 AGAACTGATG TATTATCTTG AATAAAGACT TTGTCTTGTT TACTGCCCTA

5453 AAAAAAAAA ·AAAAAAAA

25 · :

30.

10

15

20 .

BLAST Results

No BLAST result

Medline entries

35 No Medline entry

Peptide information for frame 1

ORF from 292 bp to 3801 bp; peptide length: 1170 Category: similarity to unknown protein Classification: no clue

45

50

55

40

MREKGRHDHE RTSQSHDRRH ERREDTRGKR DREKDSREER EYEQDQSSSR
51 DHRDDREPRD GRDRRDARDT RDRRELRDSR DMRDSREMRD YSRDTKESRD
101 PRDSRSTRDA HDYRDREGRD THRKEDTYPE ESRSYGRNHL REESSRTEIR
151 NESRNESRSE IRNDRMGRSR GRVPELPEKG SRGSRGSQID SHSSNSNYHD
201 SWETRSSYPE RDRYPERDNR DQARDSSFER RHGERDRRDN RERDQRPSSP
251 IRHQGRNDEL ERDERREERR VDRVDDRRDE RARERDRERE RDRERERER
301 RERDREREKE RELERERARE REREREKERD RERDRDHDD RERERERERD
351 REKERERERE ERERERERE ERERERERE ERARERDKER ERQRDWEDKD
401 KGRDDRREKR EEIREDRNPR DGHDERKSKK RYRNEGSPSP RQSPKRRREH
451 SPDSDAYNSG DDKNEKHRLL SQVVRPQESR SLSPSHLTED RQGRWKEEDR
501 KPERKESSRR YEEQELKEKV SSVDKQREQT EILESSRMRA QDIIGHHQSE
551 DRETSDRAHD ENKKKAKIQK KPIKKKKEDD VGIERGNIET TSEDGQVFSP
601 KKGQKKKSIE KKRKKSKGDS DISDEEAAQQ SKKKRGPRTP PITTKEELVE

	WO 01/98454 PCT/IB01/02050	
	L51 MCNGKNGILE DSQKKEDTAF SDWSDEDVPD RTEVTEAEHT ATATTPGSTP 701 SPLSSLLPPP PPVATATATT VPATLAATTA AAATSFSTSA ITISTSATPT 751 NTTNNTFANE DSHRKCHRTR VEKVETPHVT IEDAQHRKPM DQKRSSSLGS	
5	BOL NRSNRSHTSG RLRSPSNDSA HRSGDDQSGR KRVLHSGSRD REKTKSLEIT B5L GERKSRIDQL KRGEPSRSTS SDRQDSRSHS SRRSSPESDR QVHSRSGSFD POL SRDRLQERDR YEHDRERERE RRDTRQREWD RDADKDWPRN RDRDRLRERE	
	951 RERERDKRRD LDRERERLIS DSVERDRDRD RDRTFESSQI ESVKRCEAKL 1001 EGEHERDLES TSRDSLALDK ERMDKDLGSV QGFEDTNKSE RTESLEAGDD 1051 ESKLDDAHSL GSGAGEGYEP ISDDELDEIL AGDAEKREDQ QDEEKMPDPL	
10	1101 DVIDVDWSGL MPKHPKEPRE PGAALLKFTP GAVMLRVGIS KKLAGSELFA 1151 KVKETCQRLL EKPKGSFILL	
15	BLASTP hits	
	No BLASTP hits available	
20	Alert BLASTP hits for DKFZphtes3_12d18, frame 1	
:	No Alert BLASTP hits found	•
	Pedant information for DKFZphtes3_l2dl8, frame l	
25	Report for DKFZphtes3_12d18-1	
30	ELENGTHD 1267 EMW3 150593.45 Epid 9.22	
35	<pre>EHOMOLI TREMBL:ABO20660_1 gene: "KIAAO853"; product: "KIAAO853 protein"; Homo sapiens mRNA for KIAAO853 protein; partial cds. 0.0 EBLOCKSI BLOO422C Granins proteins</pre>	
	TBLOCKSI BLOODOF TBLOCKSI PRODODAC TBLOCKSI PRODOBPB	
40	EBLOCKS] PRODO49D EBLOCKS] PROJO63A EBLOCKS] PROD545A EBLOCKS] BLOO048 Protamine Pl proteins	
45	EBLOCKSI PFOLI4OD EBLOCKSI PFOLI4OD EKWI All_Alpha	
	EKW3 LOW_COMPLEXITY 44-32 %	
	SEQ KDRGRDFERQREKRDKPRSTSPAGQHHSPISSRHHSSSSQSGSSIQRHSPSPRRKRTPSP	
50	PRD cccchhhhhhhhccccccccccccccccccccccccc	
	SEQ SYQRTLTPPLRRSASPYPSHSLSSPQRKQSPPRHRSPMREKGRHDHERTSQSHDRRHERR	
55	PRD cccccccccccccccccccccccccccccccccccc	
	SEQ EDTRGKRDREKDSREEREYEQDQSSSRDHRDDREPRDGRDRRDARDTRDRRELRDSRDMR SEG xx.xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx	
	·	

WO 01/98454 PCT/IB01/02050 PRD DSREMRDYSRDTKESRDPRDSRSTRDAHDYRDREGRDTHRKEDTYPEESRSYGRNHLREE SEQ SEG 5 PRD **SZRTEIRNESRNESRSEIRNDRMGRSRGRVPELPEKGSRGSRGSQIDSHSSNSNYHDSWE** SEQ SEG PRD 10 TRSSYPERDRYPERDNRDQARDSSFERRHGERDRRDNRERDQRPSSPIRHQGRNDELERD SEQ ZEG PRD ERREERRVDRVDDRRDERARERDRERERDRERERERERERDREREKERELERERARERER 15 SEQ SEG PRD րերերել անական ա SEQ 20 SEG PRD RERDKERERARDMEDKDKGRDDRREKREEIREDRNPRDGHDERKSKKRYRNEGSPSPRAS SEQ SEG 25 PRD PKRRREHSPDSDAYNSGDDKNEKHRLLSQVVRPQESRSLSPSHLTEDRQGRWKEEDRKPE SEQ SEG xxxxx------PRD 30 RKESSRRYEEQELKEKVSSVDKQREQTEILESSRMRAQDIIGHHQSEDRETSDRAHDENK SEQ SEGx PRD KKAKIQKKPIKKKKEDDVGIERGNIETTSEDGQVFSPKKGQKKKSIEKKRKKSKGDSDIS 35 SEQ SEG PRD DEEAAQQSKKKRGPRTPPITTKEELVEMCNGKNGILEDSQKKEDTAFSDWSDEDVPDRTE SEQ 40 SEG xxx-----xx PRD SEQ SEG 45 PRD STSATPTNTTNNTFANEDSHRKCHRTRVEKVETPHVTIEDAQHRKPMDQKRSSSLGSNRS SEQ SEG PRD 50 NRSHTSGRLRSPSNDSAHRSGDDQSGRKRVLHSGSRDREKTKSLEITGERKSRIDQLKRG SEQ SEG PRD 55 SEQ EPSRSTSSDRQDSRSHSSRRSSPESDRQVHSRSGSFDSRDRLQERDRYEHDRERERERRD SEG PRD

	W	/O 01/98454 PCT/IB01/02050
	SEQ SEG PRD	XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
5	SEQ SEG PRD	**************************************
10	SEQ SEG PRD	· · · · · · · · · · · · · · · · · · ·
15	SEQ SEG PRD	**************************************
	SEQ SEG PRD.	• • • • • •
20	· (No	Prosite data available for DKFZphtes3_l2dl8.1)
25	(No	Pfam data available for DKFZphtes3_12d18.1)

DKFZphtes3_1417

5 group: testis derived

DKFZphtes3_1417 encodes a novel &15 amino acid protein without similarity to known proteins.

10 The mRNA is transcribed ubiquitously.
No informative BLAST results: No predictive prosite: pfam or SCOP motife.

The new protein can find application in studying the expression 15 profile of testis-specific genes.

similarity to C.elegans BD412.3

20 see also DKFZphtes3_17n3
 perhaps complete cds.

Sequenced by BMFZ

25 Locus: unknown

30

Insert length: 3522 bp
Poly A stretch at pos. 3456, polyadenylation signal at pos. 3437

	wo	01/98454				PCT/IB01/02050
	1401	GCACTGACTA	ATGAAATGTA	TTGTTTGGTT	GTGACTGTTC	AGTCCCATGA
	1451	AAAGACCCAA	ATCAGAGATG	TGAAGCTCAC	TGCTGGCTTA	AAACCAGGAC
	1501	AGGATGCCAA	TTTAACTCAG	AAGACTCACG	TGACTCTTCA	TGGACCAGAA
	1551	CTGTGTGATG	AATCCTACCC	GGCTTTACTC	ACTGACATTC	CTGTTGGAGA
5	1P0J	CTTACATCCA	GGGGAACAGC	TGGAAAAAAT	GTTGTATGTT	CGCTGTGGAA
	1651	CAGTGGGTTC	CAGAATGTTT	CTTGTATATG	TTTCTTACCT	GATAAATACA
	1701	ACCGTTGAAG	AAAAAGAAAT	TGTTTGCAAG	TGTCACAAGG	ATGAAACTGT
	1751	AACAATTGAA	ACAGTCTTTC	CATTTGATGT	TGCGGTTAAA	TTTGTTTCTA
	1901	CCAAGTTTGA	GCACCTGGAA	AGGGTTTATG	CTGACATCCC	CTTTCTGTTG
10	1851	ATGACGGACC	TCTTAAGTGC	CTCACCCTGG	GCCCTCACTA	TTGTTTCCAG
	1901	TGAGCTCCAG	CTTGCTCCAT	CCATGACCAC	AGTGGACCAG	CTCGAGTCTC
	1951	AAGTGGACAA	TGTTATCTTA	CAGACTGGAG	AGAGTGCTAG	TGAATGCTTT
	500J	TGTCTTCAAT	GCCCATCTCT	TGGAAATATT	GAAGGTGGAG	TAGCAACCGG
	2051	GCATTATATT	ATCTCTTGGA	AAAGGACCTC	AGCAATGGAG	AATATCCCCA
15	5707	TCATCACAAC	TGTCATCACT	CTGCCGCACG	TGATTGTGGA	GAATATCCCT
	5727	CTCCATGTGA	ATGCAGATCT	GCCGTCATTT	GGGCGTGTCA	GAGAGTCGTT
	5507	ACCTGTCAAG	TATCACCTAC	AGAATAAGAC	CGACTTAGTT	CAAGATGTAG
	5527	AAATTTCTGT	GGAGCCCAGT	GATGCCTTCA	TGTTCTCAGG	TCTCAAACAG
	5307	ATTCGATTAC	GTATCCTCCC	TGGCACGGAG	CAGGAAATGC	TATATAATTT
20	2351	CTATCCTCTG	ATGGCTGGAT	ACCAGCAGCT	GCCATCTCTC	AACATCAACT
	2403	TGCTTAGATT	TCCTAACTTC	ACAAATCAGC	TGCTCAGGCG	TTTTATACCT
	2451	ACCAGTATTT	TTGTCAAGCC	ACAGGGTCGA	CTCATGGATG	ATACCTCTAT
	2501	TGCTGCTGCA	TGATGTTCAA	GACCGGCCCT	TGGCTGTTGT	TACAGAGATG
0.5	2551	TTGGGCAGAG	CTATGCAGGT	GTTTCATTGT	GAACTCTAGC	TTTGATCATG
25	5601	GTAAAAAGTT	AACCTTTTCT	ATTTTTTAAT	GGATGTTATA	CCAACTATTC
	2651	AGAGGAACTC CTAATAAATA	ATACTTCAAA	AATATTAGGA	AAATCTGTCT	TATAGTTTCT
	2701 2751	GTTATTGTTG	TCTGAAATCT AAAGTCATTT	CAGTACGACA GATGAATGGT	TGAAAGAATG AAATTCTATG	TCAGACCATT
	5907	GATTTGCATG	TATAATATCA	GGAAAATTAA	GCATCCCAAG	AAAAGTAAGT TGTGACTGGA
30	2851	CAAAGAGAGC	AGATGCACCA	GTGCCTGTGC	CATAAAGTTC	CGAATCCCCC
	2901	ATGTGTCTCT	TTCAGAGCTG	GCCAGACCGG	AAATAAATCA	TTCTCATAAA
·	2951	TTCAGTGTGT	ACTCAGAACA		AACATAGGGA	GTTGTATGAC
	3007	TGATACGGAA	AACTTCCAGA		CAAAGCAGTT	TAATTAAGGT
	3051	ATCAAAAATA	TCTTTGCTTA	CTATCAAGAA	GTGTCAAATA	GGTTCAGCTT
35	3707	GCTGCCAAAA	TATGGATCAT	TTATGAAGCA	GGTTCATATT	TTAGAGGTGT
	3151	TAATAAAATC	CTCATCGGAA	AAGATCCAAA	GTGCAAGGAT	TTGATTATAA
	3507	ACATAATTTC	CTAGACTGAA	AGTTTTTGGA	AAAGATGCAG	GGTCTGAGTC
	3251	AGGCCTTCTG	GTTATATTGT	GCAGTTTCAA	AAGAACTATT	TAAAACTCTT
	3307	GAAAACTCAT	GTAAATAAAA	ATCATAGGGT	GAAAATTGTA	TTTGTTAAAA
40	3351	TACCTTAATA	ATTTAAAATG	ACCTGATTTC	CTGGAAAATT	TTATTATTCA
	3401	AAAGGTGGAG	GCATTGTAAA	AAGGAAATAG	TGATGTAAAT	AAACATGTTC
	3451	TCTTTCAAAA	AAAAAAAAA	AAAAAAAAA	AAAAAAAAA	AAAAAAAAA
	3501	AAAAAAAAA	AAAAAAAAA	AA		
45						
				BLAST Resu	ılts	
			•			

No BLAST result

50

Medline entries

No Medline entry 55

Peptide information for frame 3

ORF from 66 bp to 2510 bp; peptide length: 815 (ategory: similarity to unknown protein Classification: no clue

1 MSKQFQAFGD LFDEAIKLGL TAIQTQNPGF YYQQAAYYAQ ERKQLAKTLC 51 NHEASVMYPN PDPLETQTGV LDFYGQRSWR QGILSFDLSD PEKEKVGILA 101 IQLKERNVVH SEIIITLLSN AVAQFKKYKC PRMKSHLMVQ MGEEYYYAKD 10 . 151 YTKALKLLDY VMCDYRSEGW WTLLTSVLTT ALKCSYLMAQ LKDYITYSLE 201 LLGRASTLKD DQKSRIEKNL INVLMNESPD PEPDCDILAV KTAQKLWADR 251 ISLAGSNIFT IGVQDFVPFV QCKAKFHAPS FHVDVPVQFD IYLKADCPHP 301 IRFSKLCVSF NNGEYNGFCV IEEASKANEV LENLTGGKMC LVPGKTKKLL 351 FKFVAKTEDV GKKIEITSVD LALGNETGRC VVLNWQGGGG DAASSQEALQ 15 401 AARSFKRRPK LPDNEVHUDS IIIQASTMII SRVPNISVHL LHEPPALTNE 451 MYCLVVTVQS HEKTQIRDVK LTAGLKPGQD ANLTQKTHVT LHGPELCDES 501 YPALLTDIPV GDLHPGEQLE KMLYVRCGTV GSRMFLVYVS YLINTTVEEK : : 551 EIVCKCHKDE TVTIETVFPF DVAVKFVSTK FEHLERVYAD IPFLLMTDLL LOJ ZASPWALTIV SSELQLAPSM TTVDQLESQV DNVILQTGEZ ASECFCLQCP 651 SLGNIEGGVA TGHYIISWKR TSAMENIPII TTVITLPHVI VENIPLHVNA : · 701 DLPSFGRVRE SLPVKYHLQN KTDLVQDVEI SVEPSDAFMF SGLKQIRLRI 751 LPGTEREMLY NEYPLMAGYR RLPSLNINLL REPNETNALL RRFIPTSIEV BOJ KPQGRLMDDT SIAAA

25

BLASTP hits

30 No BLASTP hits available

Alert BLASTP hits for DKFZphtes3_1417, frame 3

No Alert BLASTP hits found

35

Pedant information for DKFZphtes3_1417, frame 3

Report for DKFZphtes3_1417.3

40

ELENGTH3 836 94249.30 EWWI 5.84 [pI]

TREMBL: CEUBO412_2 gene: "BO412.3"; Caenorhabditis **EHOMOL** elegans cosmid BO412. Le-30 [KW] Alpha_Beta EKWI LOW_COMPLEXITY 7-50 %

50

55

HIDLCKKKIGSAELSFEHDAWMSKQFQAFGDLFDEAIKLGLTAIQTQNPGFYYQQAAYYA ZEQ SEG PRD SEQ QERKQLAKTLCNHEASVMYPNPDPLETQTGVLDFYGQRSWRQGILSFDLSDPEKEKVGIL SEG PRD

	W	VO 01/98454	PCT/IB01/02050
	SEQ SEG	• • • • • • • • • • • • • • • • • • • •	KCPRMKSHLMVQMGEEYYYAKDYTKALKLLD
	PRD	հիրիկիրերի հերական անագրելու հերական հե	hhhhhhhhhhhccceeehhhhhhhhhhhh
5	SEQ SEG	YVMCDYRSEGWWTLLTSVLTTALKCSYLM	A@LKDYITYSLELLGRASTLKDD@KSRIEKN
	PRD		hhhhhhhhhhhhhhhhhhccccchhhhh
10	SEQ SEG		DRISLAGSNIFTIGVØDFVPFVØCKAKFHAP
10	PRD		hhhhhhcccceeeeeeehhhhhhhhhhcccc
	SEQ		SFNN@EYN@FCVIEEASKANEVLENLT@GKM
15	SEG PRD		ecccccceeeeecccchhhhhcccccc
	SEQ SEG		VDLALGNETGRCVVLNWQGGGGDAASSQEAL
20	PRD		eccccccceeeeeecccccchhhhhh
-0.	SEQ SEG		IISRVPNISVHLLHEPPALTNEMYCLVVTVQ:
	PRD		eeecccceeeeeccccccceeeeeeee
25	SEQ	SHEKTQIRDVKLTAGLKPGQDANLTQKTH	VTLHGPELCDESYPALLTDIPVGDLHPGEQL
	PRD		eeccccccceeeeeccccccccchh
30	SEQ		EKEIVCKCHKDETVTIETVFPFDVAVKFVST
30	PRD		eeeeeeccccceeeeeccceeeeeeh
	SEQ		IVZSELQLAPSMTTVDQLESQVDNVILQTGE
35	PRD		ehhhhhhhccceeeccccccceeeccc
	SEQ		KRTSAMENIPIITTVITLPHVIVENIPLHVN
40	PRD		ecccccceeeeeeeeeeeeccccc
	SEQ	ADLPSFGRVRESLPVKYHL@NKTDLV@DV	EISVEPSDAFMFSGLKQIRLRILPGTEQEML
	PRD		eeeccccceeeccccceeeccccccc
45	SEQ		LLRRFIPTSIFVKPQGRLMDDTSIAAA
	PRD		hhhhcccceeeeecccccccccc
50	(No	Prosite data available for DKF	Zphtes3_1417.3)
	(Ņo	Pfam data available for DKFZph	tes3_1417.3)

DKFZphtes3_15n14

5 group: testis derived

DKFZphtes3_15n14 encodes a novel 713 amino acid protein with weak similarity to the neurofilament triplet M protein of the rat-

- Neurofilaments are the intermediate filaments specific to nervous tissue. They are probably essential to the tensile strength of the neuron, as well as to transport of molecules and organelles within the axon. Until now, ESTs of the novel mRNA could only be isolated from testes, germ cells and uterus.
- No informative BLAST results: No predictive prosite: pfam or SCOP motife.

The new protein can find application in studying the expression profile of testis-specific genes.

20

similarity to neurofilament triplet M protein - rat

few EST hits (6 of 9 hits from testis)
25 perhaps complete cds.

Sequenced by GBF

Locus: unknown

30

Insert length: 2389 bp

Poly A stretch at pos. 2328, polyadenylation signal at pos. 2306

35 1 TGGGCCCCAC CTCCTCAGCA CAACTTTCTG AAAAACTGGC AGCGTAACAC 51 AGCCCTGCGG AAGAAGCAGC AGGAAGCCCT CAGCGAACAC CTAAAGAAGC LOL CAGTGAGTGA GCTGCTCATG CACACCGGGG AGACCTACAG ACGGATCCAG 151 GAGGAGCGGG AGCTCATTGA CTGCACACTT CCAACCCGGC GTGATAGGAA 201 AAGCTGGGAG AACAGTGGGT TCTGGAGTCG ACTGGAATAC TTGGGAGATG 251 AGATGACAGG TCTGGTCATG ACCAAGACAA AAACTCAGCG TGGCCTCATG 40 BOADADADA CTCATATCAG GAAGCCCCAC TCCATCCGGG TGGAGACAGG 351 ATTACCAGCC CAGAGGGACG CTTCATACCG CTACACCTGG GATCGGAGTC 401 TGTTTCTGAT CTACCGACGC AAGGAGCTGC AGAGAATCAT GGAAGAGCTG 451 GATTTCAGCC AGCAGGATAT TGATGGCCTG GAGGTGGTGG GCAAAGGGTG 501 GCCCTTCTCG GCTGTTACTG TGGAAGACTA CACAGTGTTT GAAAGAAGTC 45 551 AGGGAAGCTC CTCTGAAGAC ACAACATACT TAGGCACATT GGCCAGTTCC LOW TCTGATGTCT CCATGCCTAT TCTCGGCCCT TCTCTGCTGT TCTGTGGGAA
LOW GCCAGCTTGC TGGATCAGAG GCAGTAATCC ACAGGACAAG AGGCAGGTTG
TOW GGATTGCTGC TCACTTGACC TTTGAAACCC TAGAAGGCGA GAAAACCTCC
TSW TCAGAACTGA CTGTGGTCAA TAATGGCACC GTGGCCATTT GGTATGACTG 50 TCAGAACTGA CTGTGGTCAA TAATGGCACC GTGGCCATTT GGTATGACTG
BOL GCGACGCAG CACCAGCCGG ACACTTTCCA AGACCTTAAG AAAAACAGGA
B5L TGCAGCGATT TTACTTTGAC AACCGGGAAG GTGTGATTCT GCCTGGAGAA
TOL ATTAAAACAT TTACCTTCTT CTTCAAGTCT TTGACTGCTG GGGTCTTCAG
TSL GGAATTTTGG GAGTTTCGAA CCCATCCTAC TCTATTAGGA GGTGCTATAC
LOOL TGCAGGTCAA TCTCCACGCG GTCTCCCTGA CCCAGGACGT TTTTGAGGAT
LOSL GAGAGGAAAG TACTGGAGAG CAAGCTGACT GCCCATGAGG CAGTCACCGT
LLOL CGTTCGCGAA GTGCTGCAGG AGCTGCTGAT GGGGGTCTTG ACCCCGGAGC
LLSL GCACACCATC ACCTGTGGAT GCCTATCTCA CCGAGGAAGA CTTGTTCCGG 55

	wo	01/98454				PCT/IB01/02050
	7507	CACAGAAATC	CTCCGCTGCA	TTATGAGCAC	CAAGTGGTGC	AAAGCCTGCA
	1251	CCAACTGTGG	CGCCAGTACA	TGACCCTGCC	CGCCAAGGCT	GAGGAGGCCA
	7307	GGCCAGGGGA	CAAGGAGCAC	GTCAGCCCCA	TAGCCACAGA	GAAGGCCTCT
	1351	GTGAATGCTG		ACGCTTTAGG		CCGAAACTCA
5	1401	AGTGCCCCGG	CCTGAGAACG	AGGCCCTCAG		TCCCAGAAGG
	1451	CCAGAGTGGG	GACCAAGAGT	CCTCAGCGGA		GGAGGAGATC
	1501		AAAGCCCAGA		ACCAAGAGCC	CCTGGGAGCC
	1551	GGATGGCCTT		AGTGGAACCT		GACTTCAGAA
10	7607	AGGCAGTGAT	GGTGCTCCCT	GATGAGAACC		TGCGTTGATG
10	1651 1701	AGGCTCAACA	AAGCAGCCCT CTGCACCAGA	GGAGCTGTGC	GCTGTGGCGA	GGCCATTGCA GATGTGATTG
	1751	ACAGCCTGGT		ATGTGGCTGA		GGGCCTGCCT
	1907		CCATCTATTT			ATCAAAAATC
	1851	ACCTCCTATC	ATGGAAGTGA	AGGTACCTGT		GGGAAGGAGG
15	1901	AGCGGAAAGG		GAAAGAAGC		CAAAGACAAA
	1951	GAAGACAAGA		GCTGCTCGGG		GTCCCAACAG
	5007	CAAGAAGCAC	AAGGCAAAGG	ATGACAAGAA	AGTCATAAAA	TCTGCAAGTC
	2051	AGGACAGGTT		GACCCTACCC		CCTCTCTTCT
	5707	CAAGAACCCA		GGTCATGGGG		AGAGGCTGCA
20	2151		CGTGGGCTGC			
	5507		GCTCAGCCCC			TTTGCGCCTC
	2251		TCTCGGGCCC			
	5307	•	AAGAACCTTC			AAAAAAAA
25	2351	AAAAAAAAA		AAAAAAAAA	666666666	
23						
				BLAST Resu	ults	
30	No BL	AST result				
				Medline ent	tnice	
				HEUTING CIT		,
35						
	No Me	dline entry			•	¥-
						_
40			Peptide	information	n for frame	7
	•					
	APE f.	nom 114 bp t	to 2256 bpi	nontido los	nath: 717	•
		ory: putativ		bebrine ter	igen• ras	
45			Cell structu	re/motility	,	
73	C1033.	1,1001011	Leir Suracu	21 e7 moctito		
	1.	MHTGETYRRI	QEERELIDCT	LPTRRDRKSW	ENSGFWSRLE	YLGDEMTGLV
			MEPITHIRKP			
			LDFSQQDIDG			
50	1.51	DTTYLGTLAS	SZDVSMPILG	PSLLFCGKPA	CWIRGSNPQD	KRQVGIAAHL
			SSELTVVNNG			
	251.		EIKTFTFFFK			
			VEDANI ECAI	TAHEAUTUUR	EVIGELLMGV	1 TOFRTOCOU
	307	AVZLTQDVFE				
	301 351	DAYLTEEDLF	RHRNPPLHYE	HQVVQSLHQL	WRQYMTLPAK	AEEARPGDKE
55	301 351 401	DAYLTEEDLF HVSPIATEKA	RHRNPPLHYE SVNAELLPRF	HQVVQSLHQL RSPISETQVP	WRQYMTLPAK RPENEALRES	AEEARPGDKE GSQKARVGTK
55	301 351 401 451	DAYLTEEDLF HVSPIATEKA SP@RKSIMEE	RHRNPPLHYE SVNAELLPRF ILVEESPDVD	HQVVQSLHQL STKSPWEPDG	WRQYMTLPAK RPENEALRES LPLLEWNLCL	AEEARPGDKE GSQKARVGTK EDFRKAVMVL
55	301 351 401 451 501	DAYLTEEDLF HVSPIATEKA SPQRKSIMEE PDENHREDAL	RHRNPPLHYE SVNAELLPRF ILVEESPDVD MRLNKAALEL	HQVVQSLHQL CQKPRPLQSN	WRQYMTLPAK RPENEALRES LPLLEWNLCL LLHQMCLQLW	AEEARPGDKE GSQKARVGTK EDFRKAVMVL RDVIDSLVGH
55	301 351 401 451 501	DAYLTEEDLF HVSPIATEKA SPQRKSIMEE PDENHREDAL	RHRNPPLHYE SVNAELLPRF ILVEESPDVD	HQVVQSLHQL CQKPRPLQSN	WRQYMTLPAK RPENEALRES LPLLEWNLCL LLHQMCLQLW	AEEARPGDKE GSQKARVGTK EDFRKAVMVL RDVIDSLVGH

LOI QEKKQLGIKD KEDKKGAKLL GKEDRPNSKK HKAKDDKKVI KSASQDRFSL L51 EDPTPDIILS SQEPIDPLVM GKYTQRLHSE VRGLLDTLVT DLMVLADELS 701 PIKNVEEALR LCR

5

BLASTP hits

No BLASTP hits available

10

Alert BLASTP hits for DKFZphtes3_15n14, frame 1

No Alert BLASTP hits found

15

Pedant information for DKFZphtes3_15n14, frame 1

Report for DKFZphtes3_15n14.1

20

ELENGTHD 733

EMW1 81780.53

IpII 6.00

EBLOCKSD PF00878C

25 EBLOCKSI

BLOOL9OC DEAH-box subfamily ATP-dependent helicases

proteins

EKWI Alpha_Beta

EKWI LOW

LOW_COMPLEXITY 4-07 %

30

SEQ MHTGETYRRIQEERELIDCTLPTRRDRKSWENSGFWSRLEYLGDEMTGLVMTKTKTQRGL

35 SEQ MEPITHIRKPHSIRVETGLPAQRDASYRYTUDRSLFLIYRRKELQRIMEELDFSQQDIDG

SEG

SEQ LEVVGKGWPFSAVTVEDYTVFERSQGSSSEDTTYLGTLASSSDVSMPILGPSLLFCGKPA

SEQ CWIRGSNPQDKRQVGIAAHLTFETLEGEKTSSELTVVNNGTVAIWYDWRRQHQPDTFQDL

SEG

PRD eeeecccccchhhhhhhhhhheeeccccccceeeeecccceeeehhhhhccccchhhh

SEQ KKNRMQRFYFDNREGVILPGEIKTFTFFFKSLTAGVFREFWEFRTHPTLLGGAILQVNLH

SEG

50

45

SEQ AVSLTQDVFEDERKVLESKLTAHEAVTVVREVLQELLMGVLTPERTPSPVDAYLTEEDLF

55 SER RHRNPPLHYEHRVVRSLHRLWRRYMTLPAKAEEARPGDKEHVSPIATEKASVNAELLPRF

	•	PC1/1B01/02050
	SEQ SEG PRD	SICHERESPEKAKANING I KRITICETE AEEZADADZI KZAMEDDE
5	SEQ SEG PRD	
10	SEQ SEG PRD	······································
15	SEQ SEG PRD	WEKKQLGIKDKEDKKGAKLLGKEDRPNSKKHKAKDDKKVIKSASQDRFSLEDPTPDIILS <td< td=""></td<>
20	SEQ SEG PRD	SQEPIDPLVMGKYTQRLHSEVRGLLDTLVTDLMVLADELSPIKNVEEALRLCR
	(No	Prosite data available for DKFZphtes3_15n14.1)
25	(No	Pfam data available for DKFZphtes3_15n14.1)

5 group: cell structure and motility

DKFZphtes3_16b5 encodes a novel 268 amino acid protein with similarity to various tropomyosins.

10 Tropomyosins play regulatory roles in cellular structure and transport.

The new protein can find application in modulating cell structure and motility as well as modulationg cellular transport.

15

weak similarity to KIAAD774

perhaps complete cds.

20

Sequenced by BMFZ

Locus: unknown

25 Insert length: 1316 bp
Poly A stretch at pos. 1247, polyadenylation signal at pos. 1232

1 TGCTAAAATG GAATTAGAGA GAAGCATAGA CATCAGCAGA AGACAGAGTA 30 51 AGGAGCACAT ATGTAGAATT ACAGATCTAC AAGAGGAATT AAGACACAGA BDB GAGCATCACA TCTCTGAATT GGATAAGGAG GTTCAGCACC TTCATGAGAA 151 TATAAGTGCC CTAACCAAAG AACTGGAATT TAAGGGGAAA GAAATTCTCA 201 GAATACGAAG TGAATCTAAC CAACAGATAA GGTTGCATGA ACAAGATTTA 251 AACAAGAGAC TTGAAAAAGA GTTGGATGTC ATGACAGCAG ACCACCTCAG 35 **JDL AGAGAAAAT ATCATGCGGG CAGATTTTAA TAAGACTAAC GAGCTACTCA** 351 AGGAAATAAA TGCCGCTTTA CAAGTGTCAT TAGAAGAAAT GGAAGAAAA 401 TATCTAATGA GAGAATCAAA ACCAGAAGAT ATACAGATGA TTACAGAATT 451 AAAAGCCATG CTTACAGAAA GAGACCAGAT CATAAAGAAA CTAATTGAGG 501 ATAATAAGTT TTATCAGCTG GAATTAGTCA ATCGAGAAAC TAACTTCAAC 40 551 AAAGTGTTTA ACTCAAGTCC TACTGTTGGT GTTATTAATC CATTGGCTAA LOD GCAAAAGAAG AAGAATGATA AATCACCAAC AAACAGGTTT GTGAGTGTTC **LSI CCAATCTAAG TGCTCTGGAA TCTGGTGGAG TGGGCAATGG ACATCCTAAC** 701 CGCCTGGATC CCATTCCTAA TTCTCCAGTC CACGATATTG AGTTCAACAG 751 CAGCAAACCA CTTCCACAGC CAGTGCCACC TAAAGGGCCC AAGACATTTT 45 BD1 TGAGGTATCA GTAAGATGCA TGTGCATGAG CTCAAGGAAC ATGACTACTG 851 GAGTTTCCAT TACACATTGT TGCGTGCCTT GTAATTTTCC CCAAAGACGT 901 CCTGCTCAGA GTGAAGCTTC TCCAGTGGCT TCTCCAGATC CCCAGCGCCA 951 GGAGTGGTTT GCCCGGTACT TCACATTCTG AAAGAATTGT GTTGGCACAG DOD CTCTGTATAG ACTGTTACTA AGAGCATGAC TTTATACAGA TTGTTATGTA 50 ኔዐ5ጔ AATAGGCTTT CCTATGTCAA ACACTGTGAA TGAGAAAGTA TTTGTCTCTC 1101 CAACTTGAAA ATGCACTGTA TTTCCTGTGA TATTTATTGG AATCATTCTA 1151 TAAGGTACTA TATTATGTGT GTAATTATAA CTGTTATTTT TATTTGAGAT 1201 GGAAGAGTCT TTAACCTTTG TAATTACTGC ATAATAAATT TTGTTAGAAT 55 AAAAAA AAAAAA LOEL

BLAST Results

WO 01/98454

PCT/IB01/02050

No BLAST result

5

Medline entries

No Medline entry

10

Peptide information for frame 2

15

ORF from 8 bp to 811 bp; peptide length: 268 Category: similarity to known protein Classification: Cellular transport and traffic

20

1 MELERSIDIS RRQSKEHICR ITDLQEELRH REHHISELDK EVQHLHENIS 51 ALTKELEFKG KEILRIRSES NQQIRLHEQD LNKRLEKELD VMTADHLREK 101 NIMRADFNKT NELLKEINAA LQVSLEEMEE KYLMRESKPE DIQMITELKA 151 MLTERDQIIK KLIEDNKFYQ LELVNRETNF NKVFNSSPTV GVINPLAKQK 201 KKNDKSPTNR FVSVPNLSAL ESGGVGNGHP NRLDPIPNSP VHDIEFNSSK

251 PLPQPVPPKG PKTFLRYQ 25

BLASTP hits

30

No BLASTP hits available

Alert BLASTP hits for DKFZphtes3_16b5, frame 2

35 No Alert BLASTP hits found

Pedant information for DKFZphtes3_16b5, frame 2

40

Report for DKFZphtes3_16b5.2

ELENGTHI 270

31493.09

45 • EpII

PIR:A57013 early endosome antigen 1 - human le-O5 [HOMOL] EFUNCATI 03-19 recombination and dna repair ES. cerevisiae Y0L034wl le-05

EFUNCATI 03.22 cell cycle control and mitosis ES. cerevisiae,

50 YFR031c1 2e-05

> EFUNCATI 30.10 nuclear organization ES. cerevisiae, YFRD31c1 2e-05

EFUNCATD 11.04 dna repair (direct repair, base excision repair EFUNCATI 30.04 organization of cytoskeleton ES. cerevisiae: 55

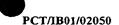
YDR356wl 7e-05

EFUNCATI 09.10 nuclear biogenesis ES. cerevisiae, YDR356w1 ?e-05



	EFUNCATE D8-07 vesicular transport (golgi network, etc.) ES. cerevisiae, YDL058wE le-04 EFUNCATE 30-03 organization of cytoplasm ES. cerevisiae,
5	YDLOSAwl le-04 [FUNCAT] I genome replication, transcription, recombination and
	repair EM. jannaschii, MJlb43D 2e-04 EFUNCATD 99 unclassified proteins ES. cerevisiae, YLR3D9cl 3e-04
10	EFUNCATI D8.16 extracellular transport ES. cerevisiae. YNL272cI 5e-04
	<pre>EFUNCATD 30.09 organization of intracellular transport vesicles</pre>
15	EKWI LOW_COMPLEXITY 4.81 % EKWI COILED_COIL 10.74 %
	SEQ AKMELERSIDISRRØSKEHICRITDLØEELRHREHHISELDKEVØHLHENISALTKELEF
20	SEGPRD cchhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
25	SEQ KGKEILRIRSESNQQIRLHEQDLNKRLEKELDVMTADHLREKNIMRADFNKTNELLKEIN SEG
	CCCCC
30	SEQ AALQVSLEEMEEKYLMRESKPEDIQMITELKAMLTERDQIIKKLIEDNKFYQLELVNRET
	PRD hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
35	SEQ NFNKVFNSSPTVGVINPLAKQKKKNDKSPTNRFVSVPNLSALESGGVGNGHPNRLDPIPN
	SEG
40	COILS
	SEQ SPVHDIEFNSSKPLPQPVPPKGPKTFLRYQ SEGxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
45	COILS
	(No Prosite data available for DKFZphtes3_1665.2)
50	(No Pfam data available for DKFZphtes3_16b5.2)
	DKFZphtes3_16p3
55	group: testis derived

WO 01/98454



DKFZphtes3_16p3 encodes a novel 1663 amino acid protein without similarity to known proteins.

The novel protein is glutamine rich and contains a cell

5 attachment RGD motif. According to the low number of ESTs and
their origin the protein seems to be expressed ubiquitously at
low levels.

No informative BLAST results; No predictive prosite; pfam or SCOP
motife.

The new protein can find application in studying the expression profile of testis-specific genes.

15 putative protein

perhaps complete cds.

Sequenced by BMFZ

20

25

10

Locus: unknown

Insert length: 5411 bp -Poly A stretch at pos. 5354, polyadenylation signal at pos. 5340

L GGCGGCCAGG TGGAGGACCT GAGCAAGCAG CTCAAGCGTG TGGACGGCCA 51 GGTGCAGGGC ATCGCCACGC ACGTGCAGCA CTTCTCCCAG GCCAGCGGGC JOB TTGACCTGGC CGCGCTAGAG TGGCCGGAGG AGCAGGAGGT GGGCGTGCGG 30 151 GCGTTCGATA GGGTGCGGAC TGGGAGTATC ATGAAGGACG CCGCCGAGGA 201 GCTCAGCTTT GCCAGGGTAC TTTTACAGCG GGTTGATGAA CTAGAGAAGC 251 TATTCAAAGA TCGGGAGCAA TTCCTGGAAC TAGTCAGCCG GAAGCTGAGT BD1 TTGGTTCCTG GTGCAGAAGA AGTCACCATG GTCACCTGGG AAGAGCTGGA 351 GCAGGCGATT ACGGACGGCT GGAGAGCCTC ACAAGCGGCC TCAGAACAC 35 403 TTATGGGATT TTCTAAGCAC GGAGGGTTCA CTTCCTTAAC ATCACCTGAA 451 GGGACTCTAA GCGGAGACTC TACCAAGCAA CCAAGTATTG AGCAGGCTCT 501 GGATTCTGCC AGTGGTCTTG GCCCGGATCG GACTGCATCA GGATCTGGTG 551 GCACAGCACA CCCCTCTGAT GGGGTTTCCA GTAGGGAACA AAGCAAGGTC LOD CCCTCTGGTA CTGGGAGACA GCAGCAGCCG AGGGCCCGTG ATGAAGCTGG 40 **653 CGTGCCACGA CTCCATCAGT CTTCTACATT CCAATTCAAA TCAGACTCAG** 701 ATCGTCACAG GAGTAGAGAG AAGCTTACCT CGACACAACC AAGAAGAAAT 751 GCACGTCCTG GTCCAGTTCA ACAGGACTTA CCCTTGGCCA GAGACCAGCC ADD CAGTAGTGTG CCCGCTAGCC AGAGTCAGGT CCATCTAAGG CCAGATCGTC **B51 GTGGGTTAGA ACCAACTGGC ATGAATCAGC CTGGATTAGT GCCTGCTAGC** 901 ACTTACCCAC ATGGTGTGGT ACCCCTCAGC ATGGGTCAGC TTGGTGTGCC 45 951 ACCACCTGAA ATGGATGATC GGGAATTGAT ACCATTTGTC GTGGATGAGC 1001 AACGTATGTT GCCACCATCA GTACCTGGCA GAGACCAGCA AGGATTGGAA 1051 CTACCTAGCA CAGACCAACA TGGTCTGGTT TCAGTCAGTG CATATCAGCA LIDI TGGTATGACA TTTCCTGGCA CAGACCAACG CAGTATGGAA CCACTTGGCA 50 1151 TGGATCAGCG TGGATGTGTA ATATCAGGCA TGGGTCAGCA AGGACTAGTA 12D1 CCCCCTGGTA TAGACCAGCA AGGATTGACA TTGCCTGTCG TCGATCAACA 1251 TGGCCTGGTT CTACCTTTTA CAGACCAGCA TGGTTTGGTA TCACCTGGTT LIBLI TGATGCCAAT TAGTGCAGAT CAGCAAGGTT TTGTGCAGCC CAGTTTGGAA 1351 GCAACTGGCT TCATACAACC TGGCACAGAG CAGCATGATT TGATCCAGTC 1401 TGGCAGATTT CAGCGTGCTT TGGTGCAGCG TGGTGCATAT CAGCCTGGCT 55 1451 TGGTCCAACC TGGTGCAGAT CAGCGTGGTT TGGTCCGGCC TGGAATGGAT 1501 CAGTCTGGTT TGGCCCAACC TGGTGCAGAT CAGCGTGGTT TGGTCTGGCC 1551 TGGAATGGAT CAGTCTGGTT TGGCCCAACC TGGTAGAGAT CAGCATGGTT

```
LLOL TGATCCAGCC TGGCACAGGT CAGCATGATT TGGTCCAATC TGGCACAGGT
LLSL CAGGGTGTCT TGGTACAGCC TGGTGTAGAT CAGCCTGGCA TGGTCCAACC
L70L TGGCAGATTT CAGCGTGCTT TGGTGCAGCC TGGTGCATAT CAGCCTGGCT
L75L TGGTCCAACC TGGTGCAGAT CAGATTGATG TGGTGCAACC TGGTGCAGAT
L80L CAGCATGGTT TGGTACAATC TGGTGCAGAT CAGAGTGATT TGGCTCAACC
L85L TGGTGCAGTT CAGCATGGTT TGGTCCAACC TGGAGTAGAT CAGCGTGGTT
L90L TGGCACAACC TCGTGCAGAT CATCAGCGTG GTTTGGTCCC ACCTGGTGCA
         1951 GATCAGCGTG GTTTGGTCCA ACCTGGTGCA GATCAGCATG GTTTGGTCCA 2001 ACCTGGAGTG GATCAGCATG GTTTGGCACA ACCTGGTGAA GTTCAGCGTA
          2051 GTTTGGTGCA ACCTGGTATA GTTCAGCGTG GTTTGGTGCA ACCTGGTGCA
 10
         2101 GTTCAGCGTG GTTTGGTGCA ACCTGGTGCA GTTCAGCGTG GTTTGGTCCA
         2151 ACCTGGAGTG GATCAGCGTG GTTTGGTTCA ACCTGGTGCA GTTCAGCGTG
         2201 GTTTGGTCCA ACCTGGTGCA GTTCAGCATG GTTTGGTCCA ACCTGGTGCA
         2251 GATCAGCGTG GTTTGGTCCA ACCTGGAGTG GATCAGCGTG GTTTGGTGCA
         23D1 ACCTGGAGTG GATCAGCGTG GTTTGGTCCA ACCTGGAATG GACCAGCGTG
15
         2351 GTTTGATCCA ACCTGGTGCA GATCAGCCTG GTTTGGTCCA GCCTGGTGCA
         2401 GGTCAGCTGG GTATGGTGCA GCCTGGAATA GGTCAGCAAG GTATGGTGCA
         2451 ACCTCAGGCA GATCCACATG GCCTGGTACA ACCTGGTGCC TATCCTCTG
         25D1 GTTTGGTACA ACCTGGTGCA TATTTGCATG ATTTATCTCA ATCTGGGACA
         2551 TATCCACGTG GTCTGGTGCA GCCAGGAATG GATCAGTATG GTTTGAGACA
20
         2603 ACCTGGTGCA TATCAGCCAG GCTTGATAGC ACCAGGCACA AAGCTTCGTG
         2651 GCTCTTCAAC ATTCCAGGCA GATTCTACAG GTTTTATATC AGTACGTCCA
         2701 TATCAACATG GTATGGTACC TCCTGGCAGA GACAATACG GCCAGGTGTC
2751 ACCACTCCTA GCCAGTCAAG GTTTGGCATC ACCTGGTATA GATCGAAGGA
2801 GTTTGGTACC ACCAGAAACT TATCAGCAAG GTTTGATGCA TCCTGGCACA
2851 GACCAGCACA GCCCAATACC ACTGAGTACA GGTTTGGGAT CTACACACCC
2901 AGATCAACAG CATGTGGCAT CACCTGGCCC AGGTGAGCAT GACCAGGTAT
2951 ACCCAGATGC AGCTCAGCAT GGCCATGCTT TCTCTCTTT TGACAGTCAT
3001 GATTCAATGT ATCCTGGTTA TCGTGGCCCA GGGTATCTAA GTGCTGATCA
25
      30
40
45
50
        4301 ACAACAAGCT GGACCGCCTG GAGCTGGACC CAGTGAAGCA GTTGCTGGAG
4351 GATCGGTGGA AATCGCTGCG ACAGCAGCTC AGGGAGCGCC CCCCACTCTA
4401 CCAGGCAGAC GAGGCGGCTG CCATGCGGAG GCAGCTCCTG GCACATTTCC
55
        4451 ACTGCCTCTC ATGTGACCGG CCCTTGGAGA CACCTGTGAC TGGACATGCC
```

WO 01/98454 PCT/IB01/02050 4501 ATCCCCGTGA CCCCCGCGGG TCCAGGCCTA CCTGGGCACC ATTCCATCCG 4551 CCCCTACACG GTGTTTGAAC TGGAGCAGGT CCGGCAGCAT AGCCGCAACC 4601 TCAAGCTGGG CAGCGCCTTC CCTCGGGGTG ACCTGGCGCA GATGGAGCAG
4651 AGCGTGGGGC GCCTGCGCTC CATGCACTCC AAGATGCTGA TGAACATTGA
4701 GAAGGTGCAG ATCCACTTCG GGGGCTCCAC CAAGGCCAGC AGCCAGATAA . 5 4751 TCCGCGAGCT GCTGCACGCC CAGTGCCTGG GCTCCCCCTG CTACAAACGG 48D1 GTGACAGATA TGGCTGATTA CACCTACTCA ACTGTGCCCC GGCGCTGCGG
4851 GGGCAGCCAC ACCCTCACCT ACCCCTACCA CCGCAGCCGC CCGCAGCACC
49D1 TTCCCCGGGG CCTGTATCCT ACTGAAGAGA TCCAGATTGC CATGAAGCAT 10 4951 GATGAGGTGG ACATCTTGGG CCTGGATGGC CACATTTACA AGGGACGGAT 5001 GGACACAAGG CTGCCAGGCA TCCTCCGAAA AGACAGCTCA GGGACCTCAA 5051 AGCGCAAGTC CCAGCAGCCC AGGCCCCACG TGCACAGGCC GCCATCCCTC 5101 AGCAGCAATG GCCAGCTGCC CTCTCGGCCA CAGAGCGCCC AGATTTCGGC 5151 TGGCAACACC TCAGAAAGAT AGACCTTCCT CCGAGGGCCG TCTCTCCCAG 5201 CCGAACACAG CCCACCCGCC CAGCTCCGCC TCGGTGGCAA ACAGGGGGCT 15 5251 GGAGAGGCAC GTGGACATGC CTCCTGGGGA GGGGCTCGAG GAGCCCACGC 5301 GGGGGCCGCG GTCCAGCACC GCTCAGTGAG CGGAGGTGTA AATAAACATT 54D3 AAAAAAAAAA A 20 BLAST Results . 25 No BLAST result Medline entries

30

No Medline entry

Peptide information for frame 1

ORF from 181 bp to 5169 bp; peptide length: 1663

Category: putative protein 40 Classification: no clue

Prosite motifs: RGD (1482-1484)

1 MKDAAEELSF ARVLLQRVDE LEKLFKDREQ FLELVSRKLS LVPGAEEVTM 51 VTWEELERAI TDGWRASRAG SETLMGFSKH GGFTSLTSPE GTLSGDSTKR 45 LOL PSIEGALDSA SGLGPDRTAS GSGGTAHPSD GVSSREGSKV PSGTGRGGGP 151 RARDEAGVPR LHQSSTFQFK SDSDRHRSRE KLTSTQPRRN ARPGPVQQDL 201 PLARDQPSSV PASQSQVHLR PDRRGLEPTG MNQPGLVPAS TYPHGVVPLS 251 MGQLGVPPPE MDDRELIPFV VDEQRMLPPS VPGRDQQGLE LPSTDQHGLV 301 SVSAY@HGMT FPGTD@RSME PLGMD@RGCV ISGMG@@GLV PPGID@@GLT 50 351 LPVVDQHGLV LPFTDQHGLV SPGLMPISAD QQGFVQPSLE ATGFIQPGTE 401 QHDLIQSGRF QRALVQRGAY QPGLVQPGAD QRGLVRPGMD QSGLAQPGAD 451 QRGLVWPGMD QSGLAQPGRD QHGLIQPGTG QHDLVQSGTG QGVLVQPGVD 501 QPGMVQPGRF QRALVQPGAY QPGLVQPGAD QIDVVQPGAD QHGLVQSGAD 551 QSDLAQPGAV QHGLVQPGVD QRGLAQPRAD HQRGLVPPGA DQRGLVQPGA 55 LOI DAHGLVAPGV DAHGLAAPGE VARSLVAPGI VARGLVAPGA VARGLVAPGA LSI VARGLVAPGV DARGLVAPGA VARGLVAPGA VAHGLVAPGA DARGLVAPGV 701 DARGLVAPGV DARGLVAPGM DARGLIAPGA DAPGLVAPGA GALGMVAPGI

			٠
			140

	·
	WO 01/98454 PCT/IB01/02050
	751 GQQGMVQPQA DPHGLVQPGA YPLGLVQPGA YLHDLSQSGT YPRGLVQPGM
	80% DQYGLRQPGA YQPGLIAPGT KLRGSSTFQA DSTGFISVRP YQHGMVPPGR 85% EQYGQVSPLL ASQGLASPGI DRRSLVPPET YQQGLMHPGT DQHSPIPLST
	901 GLGSTHPDQQ HVASPGPGEH DQVYPDAAQH GHAFSLFDSH DSMYPGYRGP
5	951 GYLSADQHGQ EGLDPNRTRA SDRHGIPAQK APGQDVTLFR SPDSVDRVLS
	1001 EGSEVSSEVL SERRNSLRRM SSSFPTAVET FHLMGELSSL YVGLKESMKD
	1051 LDEEQAGQTD LEKIQFLLAQ MVKRTIPPEL QEQLKTVKTL AKEVWQEKAK 1101 VERLQRILEG EGNQEAGKEL KAGELRLQLG VLRVTVADIE KELAELRESQ
•	1351 DRGKAAMENS VSEASLYLQD QLDKLRMIIE SMLTSSSTLL SMSMAPHKAH
10	1201 TLAPGQIDPE ATCPACSLDV SHQVSTLVRR YEQLQDMVNS LAVSRPSKKA
	1251 KLQRQDEELL GRVQSAILQV QGDCEKLNIT TSNLIEDHRQ KQKDIAMLYQ 1301 GLEKLEKEKA NREHLEMEID VKADKSALAT KVSRVQFDAT TEQLNHMMQE
	1301 GLEKLEKEKA NREHLEMEID VKADKSALAT KVSRVQFDAT TEQLNHMMQE 1351 LVAKMSGQEQ DWQKMLDRLL TEMDNKLDRL ELDPVKQLLE DRWKSLRQQL
	1401 RERPPLYCAD EAAAMRRCLL AHFHCLSCDR PLETPVTGHA IPVTPAGPGL
15	1451 PGHHZIRPYT VFELEQVRQH SRNLKLGSAF PRGDLAQMEQ SVGRLRSMHS
	1501 KMLMNIEKVQ IHFGGSTKAS SQIIRELLHA QCLGSPCYKR VTDMADYTYS 1551 TVPRRCGGSH TLTYPYHRSR PQHLPRGLYP TEEIQIAMKH DEVDILGLDG
	1601 HIYKGRMDTR LPGILRKDSS GTSKRKSQQP RPHVHRPPSL SSNGQLPSRP
	1651 QSAQISAGNT SER
20	
	BLASTP hits
25	No BLASTP hits available
	Alert BLASTP hits for DKFZphtes3_16p3, frame 1
	No Alert BLASTP hits found
30	Pedant information for DKFZphtes3_16p3, frame 1
	reduce 217 of modeloff for DRI 2price33_maps 4 frame m
	Report for DKFZphtes3_16p3.1
·35	report for Σκιζρπτες3_περ3.π
	·
	ELENGTHD 1723 EMWD 187354.98
·	EpII 6.19
40	<pre>EHOMOLI TREMBL:AFD25461_4 gene: "MO1D1.5"; Caenorhabditis</pre>
	elegans cosmid MOLDL. le-47
	<pre>EFUNCATI 30.03 organization of cytoplasm ES. cerevisiae. YDLO58wl 8e-07</pre>
	<pre>EFUNCATI O8.07 vesicular transport (golgi network, etc.)</pre>
45	cerevisiae, YDLO58wD 8e-07
	EFUNCATI 99 unclassified proteins
	[FUNCAT] 11.04 dna repair (direct repair, base excision repair
	and nucleotide excision repair) [S. cerevisiae, YKRO95w] [].[]]
50	EFUNCATD 30.10 nuclear organization ES. cerevisiae, YKRO95wD
	EBLOCKZJ PROJUGAC
	EBFOCK23 LYG3919C
	EBLOCKSI PROD543H
55	EBLOCKS PRODEIOG
	EBLOCKZI BPO453PV
	EPIRKWI RNA binding 3e-06

	wo	01/984	54						1	PCT/IB01/02	050
	EPIRK EPIRK EPIRK	(MI		hydro endop ATP 2	xylysin	e Ze reti	-10 culum 7	Pe-18			
5	EPIRK EPIRK	< W II		phosp seed	hoprote 4e-34		e-06				
	CPIRK	C LU JB		glyco	a 2e-10 protein	2e-3					
10	EPIRK EPIRK EPIRK	<wib< td=""><td></td><td>alter</td><td>otrimer native p 2e-06</td><td></td><td></td><td>e-10</td><td></td><td></td><td></td></wib<>		alter	otrimer native p 2e-06			e-10			
	EPIRK	CW JB		extra	ge prot cellula	r ma	trix 2e	-10		,	
15	EPIRK EPIRK ESUPF	<wib< td=""><td>myos</td><td>prote</td><td>ane pro in bios or doma</td><td>ynth</td><td>esis 7e</td><td>e-18</td><td></td><td></td><td></td></wib<>	myos	prote	ane pro in bios or doma	ynth	esis 7e	e-18			
	ESUPF ESUPF	EMA	glute	tin 2e enin 2 in be≥		n Ze	-NL				
20	ESUPF 3e-Di	EMA	unass	signed	ribonu	cleo	proteir	n repea	t-conta	ining p	roteins
	ESUPF ESUPF EPROS		ribo		ch prot protein			nology	3e-06		
25	EKM] EKM] EKM]		LOW_	Alpha COMPLE ED_COI	XITY		84 % 80 %				
	SEQ	GGQVI	EDLSKI	2LKRVD	GQVQGIA	THVQ	HFSQASO	SLDLAAL	EWPEEQE	VGVRAFD	RVRTGSI
30	SEG PRD COILS	ccccl								eeeeeee	
		• • • •		• • • • •	• • • • • •	• • • •	• • • • • •		• • • • • •	•••••	• • • • • • •
35	SEQ SEG PRD						• • • • • •			EVTMVTW	
	COILS		• • • • •				•••••		•••••	.cnnnnn	
40	SEQ	TDGWF	RASQAO						QPSIEQA	LDSASGL	GPDRTAS
	SEG PRD COILS		cccc		eccccc				hhhhhhh	hhhcccc	ccceeec
45				• • • • •	· · · · · · · · · · · · · · · · · · ·	• • • •	• • • • •	• • • • • •			• • • • • • •
	SEG							• • • • • •	• • • • • •	FQFKZDS	• • • • • •
50	PRD COILS	3									
	SEQ	KLTST	raprri	NARPGP	VQQDLPL	ARDQ	AAVZZA	SQSQVHL	RPDRRGL	.EPTGMNQ	PGLVPAS
55	SEG PRD	cccc								cccccc	
	COILZ		,						• • • • • •		• • • • • • •

	W	O 01/98454 PCT/IB01/0205	
	SEQ	TYPHGVVPLSMGQLGVPPPEMDDRELIPFVVDEQRMLPPSVPGRDQQGLELPSTI	Dangly
	ZEG		
	PRD	ccccccccccccccccccccccccccccccccccccccc	ccccc
	COIL	2	
5			
	SEQ	SVSAY@HGMTFPGTD@RSMEPLGMD@RGCVISGMG@@GLVPPGID@@GLTLPVV]	
	SEG	•••••••••••••••••••••••••••••••••••••••	
	PRD	_cccccccccccccccccccccccccccccccccccccc	cccc
10	COIL		
		••••••	
-	654	A DETERMINE HERE! METERS ASSETTABLE TO THE SECOND OF THE S	
	SEQ	LPFTDQHGLVSPGLMPISADQQGFVQPSLEATGFIQPGTEQHDLIQSGRFQRALV	IQRGAY
15	SEG		
15	PRD COIL:	_ccccccccccccccccccccccccccccccccccccc	cccc
	COIL.		
	SEQ	QPGLVQPGADQRGLVRPGMDQSGLAQPGADQRGLVWPGMDQSGLAQPGRDQHGL	T 4 D C T C
20	SEG	WIGEARI GYNKIGEALLEGINKZGEAKLGYNKIGEAMLGIINKZGEAKLGYNKUEF	
20	PRD		
	COILS	7	
			
25	SEQ	QHDLVQSGTGQGVLVQPGVDQPGMVQPGRFQRALVQPGAYQPGLVQPGADQIDV\	JOPGAD
	SEG	*************************************	
	PRD	ccccccccccccccccccccccccccccccccccccccc	cccc
	COILS	2	

30			
	SEQ	QHGLVQSGADQSDLAQPGAVQHGLVQPGVDQRGLAQPRADHQRGLVPPGADQRGL	
	SEG		
	PRD	_cccccccccccccccccccccccccccccccccccccc	cccc
25	COILS		
35		•••••••••••••••••••••••••••••••••••••••	• • • • •
	SEQ	DQHGLVQPGVDQHGLAQPGEVQRSLVQPGIVQRGLVQPGAVQRGLVQPGAVQRGL	
	SEG	DRUGEAREGANGUETAREGEARVZEAREGEARGEARGEARGEARGEARGEARGEARGEARGEAR	.vapcv
	PRD		
40	COILS		:cccc
-10	CVILL		
	SEQ	DARGLVAPGAVARGLVAPGAVAHGLVAPGADARGLVAPGVDARGLVAPGVDARGL	VADCM
	SEG	**************************************	
45	PRD	CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC	
	COILS	5	
	SEQ	DARGLIAPGADAPGLVAPGAGALGMVAPGIGAAGMVAPAADPHGLVAPGAYPLGL	.VQPGA
50	SEG		
	PRD	ccccccccccccccccccccccccccccccccccccccc	cccc
	COILS		
55		YLHDLSQSGTYPRGLVQPGMDQYGLRQPGAYQPGLIAPGTKLRGSSTFQADSTGF	
	SEG		• • • • •
	TCILI		

ZEQ

	W	O 01/98454	PCT/IB01/02050
	SEG PRD COIL	հերորություն ու անագրություն և հերորություն և հերորություն և հերորություններ և հերորություներ և հերորություններ և հերորություններ և հերորություններ և հերորություններ և հերորություններ և հերորություններ և հերորություներ և հերորություններ և հերոր և հերորություններ և հերորություններ	
5			
	SEQ	IPVTPAGPGLPGHHSIRPYTVFELEQVRQHSRNLKLGSAFPRGDL	AQMEQSVGRLRSMHS
	PRD.	eeeeccccccccccchhhhhhhhhhhhhhccccccccch	hhhhhhhhhhhhhh
10		•••••••••••••••••••••••••••••••••••••••	••••••
	SEQ	KMLMNIEKVQIHFGGSTKASSQIIRELLHAQCLGSPCYKRVTDMA	DYTYSTVPRRCGGSH
15	PRD COIL:		
	SEQ SEG	TLTYPYHRSRPQHLPRGLYPTEEIQIAMKHDEVDILGLDGHIYKG	• • • • • • • • • • • • • • • • • • • •
20	PRD COIL:	ccccccccccccccchhhhhhhhhcceeeecccceeee	CCCCCCCEEECCCC
		•••••••••••••••••••••••••••••••••••••••	• • • • • • • • • • • • • • • • • • • •
25	SEQ SEG	GTSKRKSQQPRPHVHRPPSLSSNGQLPSRPQSAQISAGNTSER	
23	PRD	CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC	
	COIL	zz	• • •
30			
30		Prosite for DKFZphtes3_1bp3.1	
	1002q	D16 1542->1545 RGD	PD0C0001F
35	(No F	Pfam data available for DKFZphtes3_1bp3.1)	

WO 01/98454)KFZphtes3_17ii

DKFZphtes3_17i21

5 group: transmembrane protein

DKFZphtes3_17i21 encodes a novel 224 amino acid protein without similarity to known proteins.

- 10 The novel protein contains 2 transmembrane regions. ESTs can be found in testis, retina and brain.
 No informative BLAST results; No predictive prosite, pfam or SCOP motife.
- 15 The new protein can find application in studying the expression profile of testis-specific genes and as a new marker for testicular cells.
- 20 unknown protein

Pedant: contains signal peptide(frame 1) and TRANSMEMBRANE 2 (frame 2)

25 perhaps complete cds.

Sequenced by GBF

Locus: unknown

Locus. unknown

30

Insert length: 1518 bp
Poly A stretch at pos. 1480, polyadenylation signal at pos. 1454

1 GCCAGACAGC TAGGTGTCAT TCAGGGCTGG TGTCCTCTGT CCAGGCCATC 35 51 ATGGCCTCCA CTGCCGGCTA CATCGTCTCC ACCTCCTGCA AGCACATCAT BOB TGATGACCAA CACTGGCTGT CCTCTGCCTA CACGCAATTT GCTGTGCCCT 151 ACTTCATCTA CGACATCTAC GCCATGTTCC TCTGTCACTG GCACAAGCAC 201 CAGGTCAAAG GGCATGGAGG GGACGACGGA GCGGCCAGAG CCCCGGGCAG 40 251 CACGTGGGCC ATAGCGCGTG GCTACCTGCA CAAGGAGTTC CTCATGGTGC 301 TCCACCATGC CGCCATGGTG CTGGTGTGCT TCCCACTCTC AGTGGTGTGG 351 CGACAGGGTA AGGGAGACTT CTTTCTGGGT TGCATGTTGA TGGCAGAGGT 4D1 CAGCACGCCC TTCGTCTGCC TTGGCAAGAT CCTCATCCAG TACAAGCAGC 451 AGCACACAC GCTGCACAAG GTGAACGGGG CCCTGATGCT GCTCAGCTTC 45 501 CTCTGCTGCC GGGTGCTGCT CTTTCCCTAC CTGTACTGGG CCTACGGGCG 551 CCATGCCGGC CTGCCCCTGC TGGCCGTGCC CCTGGCCATC CCTGCCCACG **LOS TCAACCTEGE CECTECECTE CTCCTEGCCC CTCAECTCTA CTEETTCTTC** L51 CTCATCTGCC GTGGGGCCTG CCGCCTCTTC TGGCCCCGCT CCCGGCCGCC 701 CCCGGCCTGC CAGGCCCAGG ACTGAGGCCG GGGGCCGGGA CCCTCCCCCT 751 CCCCACCCC ACCCCGTGG AGACAGGGCT CTGGGGCTGA TGGCTGGGGT 50 BOD TGGGAGCCAG GGTCCTCTTG CCCGGACAAC CCCAGGACTG ACGATGACCC BUL TGGGAGCCAG GGTCCTCTTG CCCGGACAAC CCCAGGACTG ACGATGACCC
B51 CGAAAGGAA GAGGCCCCAT CTCTCGGGGA CTGAGGGGGT GGAGAGAGGG
TD1 GACCTCTTCC CCCTACTCTG CCCCCTTCCT GCACACCCTT GCGCTGGAGG
T51 AGGGGAGGGG GCACCGCCTC CCACCCACTG AGGGCAGGAG GGCTTGTGGG
LD1 GAGGGACACC AACAGGGTTT CAAGGGGACC AGGAGTCAGA ATGTGGGGAG
LD51 ACGCCTCTGC CAAGGCCATC CCAGCCCCTA TGCTGCCATC CCCCAGGGCT
LLD1 CCCCATCACC CGAGAGGAGA GGACGCCCCA ACTAACCCCC GCTGGCCCTC
LL51 GGGCCTCCCG AGTGGCCGGC TGCAACCACG GCTCCTCTCC AGGGTAGGCC 55

WO 01/98454 PCT/IB01/02050 1201 AGCTTGAGGA ATCTTATTTA TTTTATTTAT TTACCCAAAT TTGAACTAGT 1251 CTGTTGGGTT GGGGGAAGGA GGTGGCTGCT ACCCCCAAGC CTTCCCAGTG 1301 CTGACAACCC CGGGGGCAGG CGAGGGCGCC CAGTCCCTCA CCATCGGCTG 1351 CACATCGCGC CCTCGGGCCC TGCCATGTCC CTGGTGCTAC TGACCTCTCA
1401 AGGCTTCCTC CAATCTGGGG TCGGGGGACC CTGGGAGGTG CTTTACAGAC 5 1451 CGCTAATAAA AGACGATCTG CGTGAACGCC AAAAAAAAA AAAAAAAAA 1501 AAAAAAAAA AAAAAAA 10 **BLAST** Results No BLAST result 15 Medline entries _____ No Medline entry 20 Peptide information for frame 3 ---------25 ORF from 51 bp to 722 bp; peptide length: 224 Category: putative protein Classification: Transmembrane proteins unclassified 30 1 MASTAGYIVS TSCKHIIDDQ HWLSSAYTQF AVPYFIYDIY AMFLCHWHKH 51 QVKGHGGDDG AARAPGSTWA IARGYLHKEF LMVLHHAAMV LVCFPLSVVW 101 RAGKADFFLA CMLMAEVSTP FVCLGKILIA YKAAHTLLHK VNGALMLLSF 151 LCCRVLLFPY LYWAYGRHAG LPLLAVPLAI PAHVNLGAAL LLAPQLYWFF 201 LICRGACRLF WPRSRPPPAC QAQD 35 BLASTP hits 40 No BLASTP hits available Alert BLASTP hits for DKFZphtes3_17i21, frame 3 No Alert BLASTP hits found 45 Pedant information for DKFZphtes3_17i21, frame 3 Report for DKFZphtes3_17i21.3 50

ELENGTHI 224

EMWD 25224-11 EpID 9-03

55 [HOMOL] TREMBLNEW:AF181646_1 gene: "BcDNA.GH12326"; product: "BcDNA.GH12326"; Drosophila melanogaster BcDNA.GH02340 (BcDNA.GH02340) mRNA, complete cds. 9e-20 [BLOCKS] PROD632H

	EBLO	CKZI	PROO5	104A			
	EBLO	CKST	BLOLE	143C			
	EKWI			MEMBR	A NIE	-	
	EKM]						
_	FKMT		rom_c	OMPLE	XIIX		6·25 %
5							
	SEQ	MASTA	IGYIVS	TZCKH.	IIDDG	SHMT:	SSAYTQFAVPYFIYDIYAMFLCHWHKHQVKGHGGDDG
	SEG			• • • • •			
	PRD	CCCCE	eeeec	ccce	eecch	ihhhl	hhhhhhheeehhhhhhhhhhhhhhhhhhhccccccc
10	MEM						***************************************
	SEQ	AARAF	AUTZD	IARGYL	HKEF	LMVI	LHHAAMVLVCFPLSVVWRQGKGDFFLGCMLMAEVSTP
	SEG				- · · · · ·		**************************************
	PRD	cccc		00000	hhhh	hhhl	hhhhhhhhhcccceeeeeccccchhhhhhhhhccc
15	MEM		·ccee		_ 1 1 7 7 1 1 1 7 M	MMM	MWWWWWWWWWWWWWWWWWWWWWWWWWWWWWWWWWWWWW
15	11611	• • • • •	• • • • •	• • • • •	[11111111	MMMMMMMMMMM
	SEQ	EUCI C	VTI TA	UVAALIT	ri i 110		ALMIL DEL CODIUL EDIN MANAGEMENT
	SEG	FVCLG	KILIK	TRUUMI	LLHK	VNG	ALMLLSFLCCRVLLFPYLYWAYGRHAGLPLLAVPLAI
			• • • • •		· • • • •	• • • •	·······································
00	PRD	cccnn	nnnnn	nnnnhr	innnc	cchi	hhhhhhhhhhheeecceeeeccc
20	MEM	• • • • •	• • • • •	• • • • •		• • • •	- MMMMMMMMMMMMMM
	ZEQ	PAHVN	LGAAL	LLAPQL	.YWFF	LICE	RGACRLFWPRSRPPPACQAQD
	SEG						•••••
	PRD	cchhh	hhhhh	hhhccc	eeee	eeco	cccccccccccc
25	MEM	• • • • •	• • • • •	• • • • • •			
	•						
	(No P	rosit	e dat	a avai	lab1	e fo	or DKFZphtes3_17i21.3)
30	(No P	fam d	ata a	vailab	le f	or I	DKFZphtes3_17i21.3)
							· · · · · · · · · · · · · · · · · · ·

DKFZphtes3_l8nl4

5 group: transcription factors

DKFZphtes3_l&nl4 encodes a novel 377 amino acid protein with similarity to human giantin.

10 Giantin is discussed as an autoantigen in rheumatoid arthritis. The novel protein contains a leucine zipper and a putative Helix-loop-helix DNA-binding domain. Therefore it might be a novel transkription factor. Most EST hits are from testis and germ cells.

The new protein can find application in modulation of gene expression and in expression profiling.

20 unknown protein

15

see DKFZphtes3<u>13</u>30i23 wrong orientation perhaps complete cds-

25
Sequenced by MediGenomix

Locus: /chromosome="16"

30 Insert length: 5282 bp
Poly A stretch at pos. 5242, polyadenylation signal at pos. 5227

1 CCGGCACCCG GAGCTCCTGG GCACACGGCA TTGGCAGGGG CCGCTTCGGC 35 51 AGAGTGATGA CTGATGATGA GTCCGAGAGC GTCCTCTCCG ACTCCCATGA 101 AGGGTCGGAG CTGGAGCTGC CTGTTATCCA GCTGTGCGGG CTGGTGGAGG 151 AGCTCAGCTA TGTAAACTCT GCTCTCAAAA CTGAGACTGA GATGTTTGAG 201 AAATATTACG CTAAACTGGA GCCCAGGGAT CAGCGACCTC CACGATTATC 251 AGAAATTAAA ATATCAGCAG CAGATTATGC ACAGTTTCGA GGCAGGCGTA 40 301 GATCCAAATC CCGGACAGGT ATGGACCGTG GGGTAGGCCT GACTGCCGAC 351 CAAAACTTG AGCTGGTACA AAAAGAGGTT GCGGACATGA AGGATGACTT 401 ACGACACAA AGGGCAAATG CGGAACGCGA CCTGCAGCAT CACGAGGCGA 451 TCATTGAGGA GGCTGAAATT CGATGGAGTG AAGTTTCGAG AGAAGTGCAT 501 GAGTTTGAAA AAGATATTCT AAAAGCCATA TCCAAGAAGA AAGGGAGTAT 45 551 TTTGGCCACT CAGAAAGTGA TGAAATACAT TGAGGACATG AACCGCCGGA **LDL GGGATAATAT GAAGGAGAAA TTACGTTTGA AAAATGTTTC TCTCAAAGTT** LSI CAGAGGAAAA AAATGCTTTT ACAATTGAGG CAGAAGGAAG AGGTGAGTGA 701 GGCCCTTCAC GATGTTGATT TTCAGCAGTT GAAGATAGAG AACGCTCAAT 751 TTCTTGAGAC AATTGAAGCA AGGAATCAAG AACTGACCCA GCTAAAGCTG BOL TCATCTGGAA ACACTCTGCA GGTTCTCAAT GCCTACAAAA GCAAGCTTCA 50 B51 CAAGGCAATG GAAATATACC TCAATCTGGA CAAGGAGATC TTGCTGAGAA PD1 AAGAGCTACT TGAAAAAATT GAAAAAGAAA CACTACAAGT AGAGGAGGAC 951 CGGGCCAAAG CCGAGGCAGT GAATAAGAGG CTCCGGAAGC AGCTGGCCGA LUDL GTTCCGGGCA CCACAGGTGA TGACTTACGT CCGGGAGAAG ATCTTAAATG 1051 CGGACCTGGA GAAGAGCATC AGGATGTGGG AAAGGAAAGT GGAGATAGCA 55 1101 GAGATGTCCT TAAAAGGCCA TCGTAAGGCT TGGAATCGAA TGAAAATAAC 1151 CAATGAGCAG TTGCAGGCAG ATTACCTTGC TGGGAAGTAG CCAGAGGCAG 1201 GCCACGGCTT ACAGACCACT ACATGACCTA TAAAAGTAAT CAGCTCCTTT

```
1251 CTAGTCACGG GCTCCTCTA CTGTTCCCTG TCTGCCTGGT GTTCCCAACC
      LEGIC CTAGGICACGG GCTCCTCTCA CTGTTCCCTG TCTGCCTGGT GTTCCCAACC
LEGIC CCCACCCAG GCTGAGTATC ATCTCCTGGG CCACATCTGC CCATGGGGAG
LEGIC TGGTTTTCACA GCCTGGCCCC TGGAACTGTT ACCACTGAAA GAACCACAGG
LHDL GCACTCTAAT GGTTTGACAC TTGTTAGCCA GCATTTAGTT CACAAGCATA
LHSL GTGAAAGTGA CCTTCCCACA CCTGGGAGAG GGATAGAGGA GGGAGAGCCA
LEGIC CCCAGTGTA TGCCATGGGC TTATCCGTGG CAGCCCCAGT GTGCAACTAT
LEGIC CACTTTACTC TTCACTTTCC TGCCGTACCT CCTGGCAC ATACTCTCCT
       LLDL CAGTTTACTC TTCAGTTTGG TGGGGTAGCT CCTGGACTAG ATACTGCTGC
LLSL AAAAGAAAAC AAGCACGAAG GAAACCAAGA TGATTTCTTC GGGCTGATAC
      1701 AACCTGTTCT GACCTGCAAA AATCCTACCT TCCCCCACCT CCCCACCGTA
1751 ATAGTCATAG TATAAGGGTT GTACAGACGC CTCAGGAGAC CTGCCTGATT
1801 CCTTTACATC CTTCTCCCTA ACATCTAGAC TATCTCTAGA GCTGTTTCCT
10
       1851 AGTCGTGAAT GCGTGATGGT CCTTCTTTGT CCCTGCAAGT ATGATCCAAC
       1901 ATGGCCCAGT TCAGAATCAG AATATGTCTT CTGTGTCATG GTGGCATTTG
       1951 GTCCATGGTG GGAGAAAGAA ATCAACTTTT CCCAGTGGTG GAGTGAGGAC
15
       2DD1 AGGGGAGGC CGGCCCTCTC AGCCTTGGAT GTGATCCATT TGCTGTAGTC
       2051 TTCCACCTTG GTGTACAGAA ACAGGCCAGG GCACGTCTCA CCACCGAAGT
       PIDL TCAGGACTCC TCTCAGAACC CACAGATCGA ACTGCTGTAG CTGGCACATC
       2151 ATTGGGCTTC CTGGGTCCCC CTGTGATAAA AGACAGAAGG CTTCAAGTCT
       2201 TAGAAAACT AGTTTTTGTT GTAAATCTAT CCTTGTGCAA TATACTGTTT
20
       2251 GTTCTAGAAA TGTTTTACGC TGGTTCTCAC TGGAAATGGG GCAAATTATA
    TABAADAAD OTTAAAADADO ADDOCCACCA CACAADAT TTAADATAG GCACCACA
       2351 ACTITICCTT TTATTATGCA AATTAGCTGT GGACTTCTGC TGATTGCCTA
       2401 TAGCTTCCTG GTTCATATTT CATTTTCTTG CCCCTTTCCA GTCCTTTGGC
25
       2451 CAAACCTTCC CTCTCTTCTG GCTTCTCATT CCTGAAATGT TGGTGTTTGT
       2501 TTCTGTTTTG TCCTGAAATG CTCACATTTT CCCTTCTCTG CCTTGCTTCA
       2551 ACCCTTAGTG TAAGCCACTT CCTGCCACCT GGCAACTGCT TACCAGCCTG
       260 GCTGGCCGTG CTCTGGGTCT TCCCTACTCC CAATGGAGCA GTCCTCTGGG
       2701 GAAGCTTGGC ATCATTAGCT AGATATGGGA CCCTGGCAAG TGACCAAATC
30
       2751 CTCTCTGAGC CAAGGTGGGA ACACAGTTAA TGCCTGTAAC ACGTGCTGAG
      2801 CACAGCACAG TGCCTGGCAC ACAGCAAACA CTCAATAGAA TATTAGCTAC
2851 CATCATCCTG ATGTCGCTAT AAAGGCCAGC ATTTTTCTGA AAAGTTGGGG
2901 AAAATGGGAA AAGCAACAAG GCAACTAGTA GGTATCACTT ACCTTACCTG
2951 CCCAGACCCC ACACCCCTAG GTTCCCCCTCC AAAGATCA CGGTACCACC
35
       BODL CATGGCCCAT CTTGGTCCGA GAAGGGGGTG GTCATCCCCA GGCTAGCCAG
      BLOL TGTTGCTCTC GCAGTTTGGA CTGAGACATG GAATGGGGCC GCAATTAACA
BLSL ACAGGAAACA ATCTGAACAG ACTGAACCAC GAGCAGCAGA AAGGCAGAAG
BEOL AGCAGCCGCT TCAGCCCCTT ACCATCCGAG ACCTGGGTGT GTGGTCTGTC
BESL TGGGCACTC ACTGGCAA TCTTCCCTCA TTCTTCTC TGTCCTAG
40
      3451 CTTGCTCTGT CCCCCATGCT GGAGTGCAGT GGCATGATCT CGGCTCGCTG
45
      BLSD CCTCATGATC CACCCGCCTC GGCCTCCCCA AGTGTTGGGA TTACAGGCGT
       3701 GAGCCACTGC ACCCGGCCTA ATTTCTGTAT TTTTAGTAGA GATGGGGTTT
50
       3751 CACGATGTTG GCCAGGCTGG TCTTAATCTA ACTTCAAGTG ATCTGCCCGC
       BBBL CTCGCCCTCT CAAAGTGCTG GGATTAGGCA TGACCA TGCCCAGTGG
       BASI GGTATTCTCT TTCAATAAAG CTCCTCTTTT CCAAGGAAGC CACACCAAAAA
       3701 CAGAGATGAA GACCAGTGGG AGACAGGG AGCAACTGGT TAGAGACC
       3951 AGCGGGGAGG CCATGCTGCA AAGCTGCCGT GATTCCCTGG TGATCTCTCA
55
       4DD1 GCAGGCCAAG GCCAGACATG TGAGGAAGGC CTTGAGGACT TCATTCTGTG
      4051 CCTCTCCTTG GATGGAAGGG GGTGCTTTAG TGTGGCACTC CTGACTTTTC
       4101 AATTGACTGG TGAAGAGGCC CTTGTGTGCA CCTCACTATG TCTGCCTAGG
```

	wo	01/98454				PCT/IB01/02050
			TCCCTGGCCA			TTCATCAGTC
			TTTGTCTTAA AGGACAAATG	CTGTAGTGGT	ATAGCCAGAG	CAAGAAAAG
		CTATGGCGTG		TCACTGTTCC	TGATTGGTAG ACCCACCTGG	GCACAGCATA
5		TACGCTTTTT			GGGGCTGCGA	CTTCTGAAGC
_			GAGGCGAACA		CCCCTCTGGA	
			AAGGGAGCGG	CCACAGCCCA		TTTCATTTTG
	4501		CCTTGACATT		CCTGACAGTG	GTAGAATAAA
10			GTGAGTGCAG CAGATTCTGA	AGTGATTCTG	CTTTTGTTGG	GTTTCAGGGA
10			AGACCTCTGG		TGATTCTACA ATATTTCCAA	TGTGGGAATT GACAGAGGAT
	4701		CGGGTCACCA	TTAAATGGTG	TGCAAGCATA	
٠:			CATGTTTAGA		AGTTAAAAAC	
•			GATTTTAAAA		TAGAGTAGAA	ATAGCTTAGA
15			GAGTCTAAGA	TACAGTTAGA	AATCAACATC	TTTGAAATTA
	4903 4953	GGGTGTGTCT	GGATTACAAA	TGATGTCAGA AAATGATGGT		
			GAGGAAGTAT		TGACGGAATG	
			AAAAGCTAT			CATTATCTAA
20			TGTTAAGATC	CCCCACCTGG	CAGAGGACCC	AGTACAAAAT
: .		AGGCACTCAA			GAAGGGCAAA	
٠.			AGTCCTTCAT AAAAAAAAAA		AAAATGACTT	TTAAAAAAAA
	2527	ааааааааа	AAAAAAAAA	AAAAAAAAA	AA	
25 -						
				BLAST Resu	ults	
	No BLA	ST result			•	
30						
	•			Medline ent	tries	
٠.	:					
35	No Med	line entry				
		-				
	•					
•			Pentido	information	n for frame	·
40				111101 1101101		
	,					
			1187 bp; p	eptide leng	th: 377	
		ry: putativ fication: r				
45			io clue .EUCINE_ZIPP	FR (19-48)		
75	110310	e mouris, c	LOCANL_ZIII			
			ZDZHEGZELE			
50			PPRLSEIKIS			
50			MKDDLRHTRA KKGSILATQK			
	ינטכ הביה	KKMLTGTEGK	EEAZEVTHDA	DEGGERTENA	GELETTEARN	GELTOI KI ZZ
			KOKI NKAMET	UI NI BUETI	DVCI I CVTCV	ETI AUCENDA
	251	GNTLQVLNAY	KZKTHKUNET	YLNLDKEILL	バンじじじじいていて	C L Q Y L L D K A
	307		KQLAEFRAPQ	VMTYVREKIL		
55	307	KAEAVNKRLR		VMTYVREKIL		

BLASTP hits

No BLASTP hits available 5 Alert BLASTP hits for DKFZphtes3_18n14, frame 3 No Alert BLASTP hits found Pedant information for DKFZphtes3_lanl4, frame 3 10 Report for DKFZphtes3_18n14-3 ELENGTHD 395 15 EMMI 46159.16 9-17 [[q] TREMBL:AF136711_1 product: "myosin heavy chain"; EHOMOLI Amoeba proteus myosin heavy chain mRNA, complete cds. 5e-06 [FUNCAT] 99 unclassified proteins ES. cerevisiae, YOR216c1 7e-04 [BLOCK2] BLOO563B Stathmin family proteins **EBFOCK2** PR00915D **IPROSITED LEUCINE_ZIPPER 1** Helix-loop-helix DNA-binding domain EPFAM3 EKW1 All_Alpha LOW_COMPLEXITY F-33 % EKWI 14.68 % EKWI COILED_COIL 30 SEQ GTRSSWAHGIGRGRFGRVMTDDESESVLSDSHEGSELELPVIQLCGLVEELSYVNSALKT _____ SEG PRD COILZ 35 ETEMFEKYYAKLEPRDQRPPRLSEIKISAADYAQFRGRRRSKSRTGMDRGVGLTADQKLE ZEQ SEG PRD 40 COILZ LVQKEVADMKDDLRHTRANAERDLQHHEAIIEEAEIRWSEVSREVHEFEKDILKAISKKK SEQ ZEG 45 PRD COILZ GSILATQKVMKYIEDMNRRRDNMKEKLRLKNVSLKVQRKKMLLQLRQKEEVSEALHDVDF SEQ 50 SEC PRD COILZ

	WO 01	/98454			PCT/IB01/02050		
	COILZ						
	• •	• • • • • • • • • • •	• • • • • • • • •				
5	SEG ELLEKIEKETLQVEEDRAKAEAVNKRLRKQLAEFRAPQVMTYVREKILNADLEKSIRMWE XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX						
10		CVEIAEMSLKGH	RKAWNRMKITI	NEQLQADYLAGK	•		
	SEG PRD hr COILS	hhhhhhhhhhhh	hhhhhhhhhh	hhhhhhhhcc	· • • •		
15							
			Prosite	for DKFZphte	s3_18n14·3		
20	P200025	37	->59 LEU(CINE_ZIPPER	P\$00007		
		•	Pfam fo	or DKFZphtes3	_18n14·3		
25	HMM_NAM	IE Helix-lo	op-helix Di	NA-binding do	nain		
30	HMM *RRRNHN	IMRERRRRndIN		PHhnVPNEKP E+ R++++ + +	_SKVEILRM ++++ +E L V+		
	++ Query LHDVDFQ	198 191 243	RRR-DNMKE	CKLRLKNVSLKVQ	RKKMLLQL-RQKEEVSEA-		
35	нмм		AIEYIrsL@x IE ++L+	k			
	Query	244	KIENAQFLE	252			

DKFZphtes3_19p12

5 group: testis derived

DKFZphtes3_19pl2 encodes a novel 664 amino acid protein without similarity to known proteins.

10 No informative BLAST results; No predictive prosite, pfam or SCOP motife.

The new protein can find application in studying the expression profile of testis-specific genes.

15

unknown protein

Sequenced by MediGenomix

20

Locus: unknown

Insert length: 2161 bp

Poly A stretch at pos. 2086, no polyadenylation signal found

```
1 CCCGAGCCAG CAACCCTGAG GGGCGGCCGG GCAGCGCCGC CACCATGTTC
            51 CTGGGCACCG GGGAGCCGGC CTTGGACACG AGTCACCTTA TCTCTCTAAG
           101 CCGAGCGTCC CTGACCCCGC AGAAGCTGTG GCTGGGAACC GCAAAGCCAG
          151 GAAGTCTGAC CCAGGCCCTG AACTCACCCC TCACCTGGGA GCATGCGTGG
201 ACTGGCGTCC CCGGCGGCAC TCCTGACTGT CTGACAGACA CCTTCAGAGT
251 GAAGAGGCCA CATCTCAGGC GCTCTGCCAG CAACGGTCAT GTCCCTGGGA
30
           BOL CTCCTGTCTA CAGAGAAAAA GAAGATATGT ATGACGAGAT TATTGAGTTA
          351 AAGAAGTCAT TGCACGTGCA GAAGAGCGAC GTGGACCTGA TGAGAACGAA
401 GCTCCGGCGC CTGGAGGAGG AAAACAGCAG GAAGGACCGG CAGATAGAGC
451 AGCTCCTGGA TCCCAGCCGC GGCACGGATT TTGTTCGGAC TCTGGCAGAG
501 AAAGGCCCG ATGCCAGTTG GGTCATAAC GGGCTGAAGC AGAGGATCCT
35
          SUL AAAAGGCCCG ATGCCAGTTG GGTCATTAAC GGGCTGAAGC AGAGGATCCT
551 GAAGCTGGAA CAGCAGTGCA AGGAGAAGGA CGGCACCATC AGCAAACTCC
601 AGACCGATAT GAAGACTACC AACCTGGAAG AGATGCGGAT CGCCATGGAG
651 ACATACTACG AGGAGGTGCA TCGTCTCCAG ACCCTCTTGG CAAGTTCTGA
701 AACCACCGGA AAGAAGCCCC TGGGGGAGAA GAAGACGGGC GCCAAAAGGC
751 AGAAGAAGAT GGGCAGTGCC CTCCTGAGCT TGTCCCGGAG TGTCCAGGAG
801 CTCACGGAAG AGAACCAGAG CCTGAAGGAG GACCTGGACC GCGTGCTGAG
851 CACCTCCCCA ACCATCTCCA AGACCAGGG TTATGTGGAG TGGAGCAAGC
40
          901 CCCGGCTGCT GAGGCGCATT GTGGAGCTGG AGAAGAAACT AAGTGTGATG
45
          951 GAGAGCTCAA AATCACACGC CGCAGAGCCA GTCAGATCAC ACCCGCCAGC
        DOD CTGCCTTGCA TCCAGCTCTG CGCTGCACAG ACAGCCACGA GGGGACCGCA
DOSD ACAAGGACCA CGAGCGTCTC CGAGGGGCTG TGAGAGACCT GAAGGAAGAG
DDD CGGACCGCGC TGCAGGAGCA GCTGCTGCAG AGAGATTTGG AGGTGAAGCA
50
         1151 GCTCCTGCAG GCGAAGGCCG ACCTGGAGAA GGAGCTGGAG TGCGCGAGGG
         1201 AGGGCGAGGA GGAGAGGAGA GAGCGAGAGG AGGTTTTGAG AGAGGAGATT
        1251 CAGACACTTA CCAGCAAGCT CCAAGAATTG CAAGAAATGA AGAAAGAAGA
        LEGA GAAAGAGGAT TGCCCGGAAG TTCCTCATAA GGCCCAAGAG CTCCCAGCTC
        1351 CCACTCCCAG CAGCAGGCAC TGCGAGCAAG ACTGGCCGCC GGATTCCAGC
55
        1401 GAGGAGGGC TCCCGCGGCC CCGCTCCCCC TGCTCTGATG GGAGAAGAGA
        1451 CGCCGCGGCC AGAGTCCTGC AGGCCCAGTG GAAGGTGTAC AAGCACAAGA
        1501 AAAAAAAGGC TGTTCTGGAT GAGGCGGCTG TGGTGCTTCA GGCAGCTTTC
        1551 AGGGGACATC TCACGCGGAC AAAGCTCTTA GCAAGCAAAG CACATGGCTC
```

	WO 01/98454		PCT/IB01/02050
	1601 AGAGCCACCC AGCGTGCCAG		
	1651 GCGTTCCGAG CCCCATCGCC		
	1701 GCCATCGTCA TCATCCAGTC 1751 GCACAGTGCT ACCGGTAAAA		
5	1801 GATCGGCTTC AGCCACACAC		
	1851 GCTCTTCCTG ACCCCTCTCC	CTCAGGGCCA CAGGCCTTGG	CACCTCTACC
		ATGATTCCGA CGATATTGTC	
	1951 CTCTGCCCAC GAAGAACTTT 2001 CCGTGATGGC AGCGCTGCCG		
10	2051 TTTATCGTGT TAGGAGAAGA		
	AAAAAAAA AAAAAAAA		
	STEP AVVVVVVV		•
15		BLAST Results	
	No BLAST result		
20			
		Medline entries	
			•
25	No Medline entry		
23			
	Peptide	information for frame	3
30			
	ORF from 45 bp to 1976 bp;	peptide length: 644	
	Category: similarity to unki Classification: unclassified	nown protein .	
	Prosite motifs: RGD (332-33		
35			
	I MELCTCEDAL DTSULTSUSE	ASI TOAKI W. ATAKOASI TA	AL MODE THE
	<pre>1 MFLGTGEPAL DTSHLISLSR 51 AWTGVPGGTP DCLTDTFRVK</pre>	RPHLRRSASN GHVPGTPVYR	
	JOJ ELKKSLHVQK SDVDLMRTKL	RRLEEENSRK DRQIEQLLDP	
40		LEQQCKEKDG TISKLQTDMK	TTNLEEMRIA
	201 METYYEEVHR L@TLLASSET 251 @ELTEEN@SL KEDLDRVLST	TGKKPLGEKK TGAKRQKKMG SPTISKTQGY VEWSKPRLLR	
	301 WESSKSHAA EPVENKVEST	LASSSALHRQ PRGDRNKDHE	
	351 EERTALQEQL LQRDLEVKQL	LQAKADLEKE LECAREGEEE	RREREEVLRE
45	407 EIGLTLZKTG ETGEWKKEEK		
	451 SSEEGLPRPR SPCSDGRRDA 501 AFRGHLTRTK LLASKAHGSE	AARVLQAQWK VYKHKKKKAV PPSVPGLPDQ SSPVPRVPSP	
	551 EEAIVIIQSA LRAHLARARH		
	POT LAALDDAZAZ CAGATABTA		
50		·	
	•	RI ASTP hite	

BLASTP hits

55 No BLASTP hits available

Alert BLASTP hits for DKFZphtes3_19p12, frame 3

Pedant information for DKFZphtes3_19pl2, frame 3

5

Report for DKFZphtes3_19pl2.3

```
ELENGTHI 644
10
    71810-41
    [pI]
              8.80
    EHOMOLI
                   TREMBL: ABD28946_1 gene: "KIAA1023"; product:
    "KIAAl023 protein"; Homo sapiens mRNA for KIAAl023 protein;
    partial cds. [].
    EFUNCATI
             30.03 organization of cytoplasm
                                               ES. cerevisiae,
    YDL058wl 2e-07
    EFUNCATI D8-D7 vesicular transport (golgi network, etc.)
    cerevisiae, YDLO58wl 2e-07
    EFUNCATD 99 unclassified proteins
                                           ES. cerevisiae YLR309cl
20
    3e-06
    EFUNCATI 30-04 organization of cytoskeleton
                                                     ES. cerevisiae,
    YDR356w3 2e-05
    EFUNCATI 09-10 nuclear biogenesis
                                           ES- cerevisiae YDR356wl
    2e-05
25
    [FUNCAT] D3.22 cell cycle control and mitosis [IS. cerevisiae.
    YDR356w3 2e-05
    EFUNCATI
             98 classification not yet clear-cut
                                                   ES. cerevisiae,
    YJR134c1 4e-05
    IBLOCKS3 DMO13541
    EBLOCKS3 BLOOL27B GHMP kinases ATP-binding domain proteins
30
    IBLOCKSD BLOD326C Tropomyosins proteins
    IBLOCKSI BLOILLOB Kinesin light chain repeat proteins
    EBLOCKSI BLOOBEOD Glucoamylase proteins region proteins
    EBLOCKSI BPD4473C
35
    EBLOCKSI BLOO412B Neuromodulin (GAP-43) proteins
             3.6.1.32 Myosin ATPase 3e-08
    [EC]
    EPIRKU
                  tandem repeat 3e-08
                  transmembrane protein 2e-07
    [PIRKW]
   : [PIRKW]
                  muscle contraction 3e-O8
40 EPIRKWI
                  actin binding 3e-08
    [PIRKW]
                  ATP 3e-08
                 thick filament 3e-O8
    IPIRKUJ
    [PIRKW]
                  alternative splicing 7e-07
                  coiled coil 3e-DA
    EPIRKUI
45
    EPIRKU
                  P-loop 3e-D8
    [PIRKW]
                  heptad repeat 2e-07
                  methylated amino acid 3e-08
    EPIRKU
    [PIRKW]
                  hydrolase 3e-08
    EPIRKU
                  Golgi apparatus 2e-07.
50
    ESUPFAMD myosin heavy chain 3e-08
   ESUPFAMB
             myosin motor domain homology 3e-DB
             alpha-actinin actin-binding domain homology &e-Db
   ESUPFAMD
   ESUPFAMI
             plectin &e-Ob
   ESUPFAMD
             ribosomal protein SLD homology &e-Db
   ESUPFAMI
             giantin 2e-07
   EPROSITED RGD L
   EKW1
             All_Alpha
   [KW]
             LOW_COMPLEXITY
                               14.60 %
```

EKMJ COILED_COIL 72.55 %

5	SEQ SEG	MFLGTGEPALDTSHLISLSRASLTPQKLWLGTAKPGSLTQALNSPLTWEHAWTGVPGGTP
	PRD COIL	
10	SEG	DCLTDTFRVKRPHLRRSASNGHVPGTPVYREKEDMYDEIIELKKSLHVQKSDVDLMRTKL
	PRD COIL	ccccchhhhhhhhhhhhcccchhhhhhhhhhhhhhhhh
15		
	SEG SEQ	RRLEEENSRKDRQIEQLLDPSRGTDFVRTLAEKRPDASWVINGLKQRILKLEQQCKEKDG
20	PRD COIL	
20		ccccc
	SEQ SEG	TISKLØTDMKTTNLEEMRIAMETYYEEVHRLØTLLASSETTGKKPLGEKKTGAKRØKKMG
25	PRD COIL	
		cccc
	SEQ	SALLSLSRSVQELTEENQSLKEDLDRVLSTSPTISKTQGYVEWSKPRLLRRIVELEKKLS
30	PRD COIL	
35	SEQ SEG	VMESSKSHAAEPVRSHPPACLASSSALHRQPRGDRNKDHERLRGAVRDLKEERTALQEQL
	PRD COIL:	
	4	•••••••••••••••••••••••••••••••••••••••
40	SEQ	LQRDLEVKQLLQAKADLEKELECAREGEEERREREEVLREEIQTLTSKLQELQEMKKEEK
	PRD COIL:	
45		
	SEQ	EDCPEVPHKAQELPAPTPSSRHCEQDWPPDSSEEGLPRPRSPCSDGRRDAAARVLQAQWK xx
	PRD COILS	
50		ccccc
	SEQ	VYKHKKKAVLDEAAVVLQAAFRGHLTRTKLLASKAHGSEPPSVPGLPDQSSPVPRVPSP
55	PRD COILS	hhhhhhhhhhhhhhhhhhhhcchhhhhhhhhhccccccc
	SEQ	IAQATGSPVQEEAIVIIQSALRAHLARARHSATGKRTTTAASTRRSASATHGDASSPPF

	WO 01/98454	PCT/IB01/02050		
	SEGxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx	xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx		
5	*********************************	• • • • • • • • • • • • • • • • • • • •		
,	SEG LAALPDPSPSGPQALAPLPGDDVNSDDSDDIVIAPSLP			
	PRD eeeccccccccccccccccccccceeeeecccc	cccc		
10	COILS	• • • • • • • • • •		
	Prosite for DKFZphtes3	_l9pl2·3		
15	PS00016 332->335 RGD	PD0CDDD1P		
	(No Pfam data available for DKFZphtes3_19pl	2•3)		

PCT/IB01/02050

DKFZphtes3_20hl2

5 group: transmembrane protein

DKFZphtes3_20hl2 encodes a novel 1204 amino acid protein without similarity to known proteins.

- The novel protein contains 1 transmembrane region and two leucine zippers.
 No informative BLAST results; No predictive prosite, pfam or SCOP motife.
- 15 The new protein can find application in studying the expression profile of testis-specific genes and as a new marker for testicular cells.
- 20 putative protein

perhaps complete cds.
Pedant: TRANSMEMBRANE]

25 Sequenced by MediGenomix

Locus: unknown

Insert length: 5894 bp

30 Poly A stretch at pos. 5874, no polyadenylation signal found

L CTCTGCCTTT CCTCTCGCAG CCACCCTTCC TCTCAGACCA GTACGGTGGC 51 CGACGGGAGT CAGACGCTGG GGATGAATGA AGGATCAACA AACAGTAATA 35 101 ATGACTGAAT GTACAAGTCT TCAGTTTGTC AGCCCTTTTG CTTTTGAGGC 151 AATGCAGAAG GTGGATGTTG TTTGCCTGGC ATCTTTAAGT GATCCAGAAT 201 TAAGACTTCT TCTGCCCTGT TTGGTACGGA TGGCACTTTG TGCACCTGCT 251 GACCAGAGCC AAAGCTGGGC TCAGGATAAG AAACTCATCC TTCGCCTTCT BOL TTCTGGAGTG GAAGCTGTCA ACTCCATTGT TGCATTGTTG TCCGTGGACT 40 351 TTCATGCTTT AGAACAAGAT GCCAGCAAAG AACAGCAGCT TAGGCATAAA 4D1 CTTGGAGGAG GCAGTGGAGA GAGCATCCTG GTATCACAGC TTCAGCATGG 451 ACTGACGTTA GAGTTTGAAC ACAGTGATTC ACCTCGTCGA TTGCGTCTTG 501 TGCTTAGTGA ACTGTTGGCA ATTATGAACA AGGTGTCTGA GTCCAACGGA 551 GAATTTTTT TCAAGTCTTC TGAACTTTTT GAGAGTCCAG TATATTTGGA 45 LOL GGAAGCTGCA GATGTACTTT GTATTTTACA AGCAGAGCTC CCTTCCTTGC 651 TCCCTATAGT TGATGTAGCT GAAGCTTTGC TACATGTTAG AAATGGTGCC 7DL TGGTTCTTGT GTCTCTTGGT GGCCAATGTT CCTGATAGTT TTAATGAAGT 751 TTGTAGGGGC CTGATAAAAA ATGGAGAACG ACAAGATGAA GAAAGTCTTG BD1 GAGGAAGGCG CAGGACAGAT GCCTTACGCT TCTTGTGTAA AATGAATCCT 851 TCTCAGGCCC TCAAGGTCCG AGGCATGGTG GTGGAAGAAT GTCACTTGCC 50 9DL AGGCCTTGGT GTGGCTTTGA CATTGGATCA TACTAAAAT GAAGCTTGTG 951 AGGATGGAGT GAGTGACTTG GTTTGTTTTG TAAGTGGTTT GCTTCTTGGA BOOL ACAAATGCGA AAGTCCGGAC TTGGTTTGGA ACTTTTATCC GAAATGGACA 1051 GCAGAGAAAA AGAGAGACCA GCAGTTCTGT CCTTTGGCAG ATGAGAAGGC 55 LIDI AGCTTCTTCT GGAGTTGATG GGCATTCTTC CCACAGTAAG AAGCACCCGA 1151 ATTGTGGAAG AAGCTGATGT GGATATGGAG CCCAATGTGT CTGTGTATTC 1201 GGGGCTGAAA GAAGAGCATG TTGTGAAAGC CAGTGCACTC TTACGTCTGT 1251 ACTGTGCTTT GATGGGGATC GCTGGACTCA AACCAACTGA AGAAGAAGCT

LBDL GAGCAATTAC TGCAGTTGAT GACGAGCCGT CCTCCTGCTA CGCCAGCTGG 1351 GGTTCGCTTT GTTTCACTTT CCTTTTGTAT GCTACTGGCC TTTTCTACAC BUDD TTGTCAGTAC ACCTGAACAG GAGCAGCTGA TGGTGGTGTG GCTAAGTTGG 1451 ATGATAAAAG AAGAAGCGTA TTTTGAGAGT ACTTCAGGCG TCTCTGCTTC 1501 TTTTGGGGAG ATGTTATTAT TGGTGGCTAT GTACTTTCAC AGCAACCAGC 5 1551 TTAGTGCTAT CATTGACTTG GTCTGTTCCA CTTTGGGGAT GAAGATTGTA JLOJ ATTAAGCCAA GCTCCTTGAG CAGGATGAAG ACAATCTTCA CACAGGAAAT 1651 TTTTACTGAG CAGGTTGTCA CAGCTCATGC AGTTCGGGTC CCTGTCACCA
1701 GCAACCTGAG TGCCAACATT ACTGGATTTT TGCCTATTCA TTGTATTTAC 1701 GCAACCTGAG TGCCAACATT ACTGGATTTT TGCCTATTCA TTGTATTTAC
1751 CAGCTTCTCA GGAGCCGTTC CTTTACCAAG CACAAAGTGT CAATAAAAGA
1801 TTGGATTTAT AGACAGCTGT GTGAAACCTC TACTCCACTT CATCCTCAAT
1851 TACTTCCTTT GATTGATGTG TACATAAATT CTATACTTAC TCCTGCGTCG
1901 AAATCTAATC CAGAAGCCAC AAATCAGCCA GTCACAGAAC AGGAGATACT
1951 CAATATTTTC CAAGGAGTCA TTGGGGGTGA CAACATCCGC CTTAATCAGC
2001 GTTTCAGTAT CACAGCACAG CTTTTGGTGC TCTACTATAT ACTGTCTTAT
2051 GAAGAGGCTC TTCTAGCAAA CACGAAGACT TTAGCTGCCA TGCAAAGAAA
2101 GCCCAAATCA TATTCTTCTT CTTTAATGGA TCAGATTCCT ATCAAATTCC
2151 TTATTCGACA GGCTCAAGGG CTGCAGCAGG AGTTGGGAGG GTTGCATTCA 10 15 2151 TTATTCGACA GGCTCAAGGG CTGCAGCAGG AGTTGGGAGG GTTGCATTCA 2201 GCTTTACTAC GTCTCCTTGC TACTAACTAC CCACATTTAT GTATTGTGGA TGTGAAGAAG AAATCACAGG GACTGATGCC CTGCTACGGC

BBL GAATGCTCCT GACTAATAAT GCTAAAAATC ATTCTCCCAA ACAACTCCAA

BBL GAAGCATTTT CAGCTGTCCC AGTAAATCAC ACACAAGTGA TGCAGATTAT

BHL AGAACACTTG ACTCTACTCT CTGCCAGTGA ACTTATACCA TATGCGGAAG

BHSL TGTTAACATC CAATATGAGC CAGCTATTGA ATTCAGGGGT TCCACGGAGA 20 TGTTAACATC CAATATGAGC CAGCTATTGA ATTCAGGGGT TCCACGGAGA

2501 ATTCTGCAAA CAGTCAATAA ACTATGGATG GTTCTTAATA CTGTGATGCC

2551 TAGAAGGCTA TGGGTAATGA CGGTTAATGC ACTTCAGCCT TCAATAAAGT

2601 TTGTACGACA ACAAAAGTAT ACTCAGAATG ACCTGATGAT AGATCCTCTC

2651 ATTGTCCTAA GGTGTGATCA GAGGGTTCAC AGATGCCCCC CACTGATGGA

2701 TATTACCCTA CACATGTTGA ATGGATATCT TCTTGCATCT AAAGCCTACC

2751 TTAGTGCTCA TCTGAAGGAA ACAGAGCAAG ATAGGCCTTC CCAGAATAAT

2801 ACAATTGGTT TAGTTGGACA AACTGATGCT CCGGAAGTTA CCAGGGAAGA

2851 ATTGAAAAAT GCATTACTGG CCGCTCAGGA TAGTGCAGCT GTCCAGATTC

2901 TCTTAGAGAT TTGCCTACCT ACTGAAGAGG AGAAAGCAAA TGGTGTCAAT

2911 AAATAAGGGA ATGGAGGAAG GAGAAGACAA TTTTGCTCTGT AACCTTCGAG 25 30 BDDL AAATAAGGGA ATGGAGGAAG GAGAAGACAA TTTGCTCTGT AACCTTCGAG
BDSL AAGTTCAGTG CCTTATCTGT TGTCTCTTGC ACCAAATGTA CATTGCAGAT 35 BIDI CCCAACATTG CTAAGCTTGT TCACTTTCAG GGTTATCCAT GTGAACTTTT 3151 GCCTCTGACG GTCGCAGGTA TTCCATCTAT GCACATCTGT CTAGATTTCA 3201 TACCTGAGCT TATTGCACAG CCAGAACTTG AGAAACAGAT ATTTGCTATC 3251 CAGTTGCTTT CTCACTTGTG TATACAATAT GCATTACCAA AGTCACTTAG 40 AATTTTAAAA ATGTCATGG AACTTTGTTA ACAGTTTTAAA 3351 CACAGGCTAA GCGGTATGCT TTTTTTATGC CAACTCTGCC AAGTTTGGTC 3401 TCTTTTTGTC GAGCATTTCC TCCATTGTAT GAGGATATTA TGTCTTTGCT 3451 GATCCAAATA GGGCAAGTTT GTGCCTCTGA TGTTGCCACT CAGACAAGAG 3501 ACATTGATCC AATTATTACA CGTCTTCAAC AAATAAAGGA GAAACCAAGT 3551 GGATGGTCTC AAATCTGTAA AGATTCATCT TATAAAAGG GATCCAGGA 45 BLOD CACTGGAAGC ATGGATCCTG ATGTACAGCT CTGTCACTGT ATTGAAAGAA 3651 CAGTAATTGA AATAATAAAT ATGAGTGTTA GTGGAATTTA AAACAAAATT 50 55 4101 TGCAAAGATG GGAGAGGAAA AAAGGGTAAA GGGAAAGGAG AATTAAGGAA 4151 ATAATAGGAG TTAAAAACAC AAGTAGAAAT CTCAAAGATT TGCAGTGCAA

WO 01/98454 PCT/IB01/02050 4201 GTAATAGTAA TGCAAGTTGG AATTCTAGTT CTCAAGAAAG AGTATTGAGA 4251 AGACTTTTAA AAAGGCAAGT AGCTTTTGTA AATGATTTCT GTGGAAATAC 4301 AGATGAGGAT TTAAAGATTT CACATATTTG CTTCAATTTT TATTAATATA 4351 TGAAGCCATA TGTTTAAAGA GATACTTGAA TAATTTGGAA TTTTAAGATA 4401 CTGGTGTAAA AGTGTTTACA GAAACATCTT TGTTCAAAGA AGAACCTGAG 5. 445% AGATCTCATT TAGTTTTATG TTTTAAATTT ATTTTTATAA TGCTTTATTA 45DL ACTTACCTAA TGCTCAGAGG GGGGAAATAT GTATCAAATT AAATGAAGGT 4551 AGAGCAATAA AACCCACTGG ATTAAAGAGC TCTTGGTTTG TCATCAGGAT 4603 TATAATTCAT ATCTTACTTT GAGAAGATCT TTGAGTAAGA AAATGCAGTG 465% TTTGAACCTG AGGAAAAGTT AAAGTGTAGA AAATATTGTC TTGCCGAAGG 10 4701 ATTTTGCAGT CCTCTGTCAG TAACTTCCAT TGATTAGGCA GACATATTCA 4751 GGTAAACCCT AATCATTAAA AAAAAATTAT CAATGTAGAA AGTAATTCCC 48D3 TTTTTTCTCT CTGAGATATA CCTCAATCAC ACACTTCCCC ACCCCCACTT 4851 GAAACAGACC TCTTCACTTG TGTTTTTTTT TTTTTTTCC TGAGGTGGAG 4901 TCTTCCCCTG TTGCCCAGGC TGGAGTGCAG TGGGATGATC TTGGCTCACT . 15 4951 GCAACTTCTG CCACCTGGGT TCAAGGGATT CTCGTGCCTC AACCTCCTGA 5001 GTAGCTGGGA CTGCAGGCAC GCGCCACCTG TATTTTTGTA TTTTTAGTAG 5051 AGACGGGGGT TTGCCATGTT GCCCAGACTG GTTTTGAACT CCTGGCCTCA 5101 GGTGATCTGC CCACCTTGGC CTCCCAAAGT GCTGGGATTA CAGGTGTGAG 5151 CCACCGCACC TGGCCAGACC GCTTCACTTG TAAAAGAAAT TAGGCTAATA 20 5201 AGAAGGTGTA GTTTTTGAGA AATGAAATTT AACTTTAGCC TTTTCACTAG 5253 TAAATAGTCA CATCTCATTT TCTTCCTTIG TAAAATGGGG TTACTACTGG 5301 CCCTACCTCA TATTCTATGA GAATGAGTTT GTAGCTGTTT CAAATCATGA 5351 AGTGCATAGT ATCACATGTG ATAGAATATT TATAACTTTT TATTAGATGC 25 - 5401 TTAATGTTCA ATTAAGTAAT TTTGATGTGA AAAATAAAAG TAATAAAAGT 5451 ATCTTAAAAA TAGCATAAGA ATTTTCATAT TTTTAAACAA GGCAGTTTTG 5501 TAGTCCCTTA AGATTAAATA CAACTGCTCC TTTTTTTTT AAACTGAGGC 5551 CTTGCGATAT TTTGTGTGAA TAGATATGCC CTAGGAGTTC AGAAAAAGTT 5601 AAAAGTATGT TTTCTAATTA AATGCAGTGC ACATTCCTGG ATCAATATTC 565 AAAGACTGGT CATAACCTGC TGTGTTAAAA TAATCACATA TGCTCTTTTT . 30 5701 CATCAGATTT GTTGATGATG TAAATAAAAT GTGTAAATAT ATTAGTAAAT 5751 GTTAATATTC ATGTATTTTA AGTTAAGGTT ATAAAATTTG TCACAATGTG 5803 TTTTTTATT CAAGTGAAAA CAGATGTGTG CAGCTATTTT GAATATTGGT

BLAST Results

40 No BLAST result

·35 ·

Medline entries

No Medline entry

50 Peptide information for frame 2

ORF from 77 bp to 3688 bp; peptide length: 1204 Category: putative protein

55 Classification: unclassified

Prosite motifs: LEUCINE_ZIPPER (167-184)

LEUCINE_ZIPPER (692-709)

```
I MKDQQTVIMT ECTSLQFVSP FAFEAMQKVD VVCLASLSDP ELRLLLPCLV
             51 RMALCAPADQ SQSWAQDKKL ILRLLSGVEA VNSIVALLSV DFHALEQDAS
           101 KERRLRHKLG GGSGESILVS RLRHGLTLEF EHSDSPRRLR LVLSELLAIM
           151 NKVSESNGEF FFKSSELFES PVYLEEAADV LCILQAELPS LLPIVDVAEA
  5
           201 LLHVRNGAWF LCLLVANVPD SFNEVCRGLI KNGERQDEES LGGRRRTDAL
           251 RFLCKMNPSQ ALKVRGMVVE ECHLPGLGVA LTLDHTKNEA CEDGVSDLVC
           301 FVZGLLLGTN AKVRTWFGTF IRNGQQRKRE TSZSVLWQMR RQLLLELMGI
351 LPTVRSTRIV EEADVDMEVN VSVYSGLKEE HVVKASALLR LYCALMGIAG
           401 LKPTEEEAER LLRLMTSRPP ATPAGVRFVS LSFCMLLAFS TLVSTPERER 451 LMVVWLSWMI KEEAYFESTS GVSASFGEML LLVAMYFHSN RLSAIIDLVC
10
           501 STLGMKIVIK PSSLSRMKTI FT@EIFTE@V VTAHAVRVPV TSNLSANITG
           551 FLPIHCIYAL LRSRSFTKHK VSIKDWIYRA LCETSTPLHP ALLPLIDVYI
551 FLPIHCIYQL LRSRSFTKHK VSIKDWIYRQ LCETSTPLHP QLLPLIDVYI
601 NSILTPASKS NPEATNQPVT EQEILNIFQG VIGGDNIRLN QRFSITAQLL
15 651 VLYYILSYEE ALLANTKTLA AMQRKPKSYS SSLMDQIPIK FLIRQAQGLQ
701 QELGGLHSAL LRLLATNYPH LCIVDDWICE EEITGTDALL RRMLLTNNAK
751 NHSPKQLQEA FSAVPVNHTQ VMQIIEHLTL LSASELIPYA EVLTSNMSQL
801 LNSGVPRRIL QTVNKLWMVL NTVMPRRLWV MTVNALQPSI KFVRQQKYTQ
851 NDLMIDPLIV LRCDQRVHRC PPLMDITLHM LNGYLLASKA YLSAHLKETE
20 901 QDRPSQNNTI GLVGQTDAPE VTREELKNAL LAAQDSAAVQ ILLEICLPTE
951 EEKANGVNPD SLLRNVQSVI TTSAPNKGME EGEDNLLCNL REVQCLICCL
1001 LHQMYIADPN IAKLVHFQGY PCELLPLTVA GIPSMHICLD FIPELIAQPE
1051 LEKQIFAIQL LSHLCIQYAL PKSLSVARLA VNVMGTLLTV LTQAKRYAFF
1101 MPTLPSLVSF CRAFPPLYED IMSLLIQIGQ VCASDVATQT RDIDPIITRL
25 1151 QQIKEKPSGW SQICKDSSYK NGSRDTGSMD PDVQLCHCIE RTVIEIINMS
1201 VSGI
         7507 AZ61
30
                                                          BLASTP hits
       No BLASTP hits available
                            Alert BLASTP hits for DKFZphtes3_20hl2, frame 2
35
       No Alert BLASTP hits found
                           Pedant information for DKFZphtes3_20hl2, frame 2
                            40
                                          Report for DKFZphtes3_20hl2.2
       ELENGTHD 1204
45
                        134347.53
       5.75
       [[q]
                                TREMBL:CEZC376_3 gene: "ZC376.6"; Caenorhabditis
       EHOMOLI
       elegans cosmid ZC376 2e-22
       EPROSITED LEUCINE_ZIPPER 2
50
       TRANSMEMBRANE 1
                       LOW_COMPLEXITY
       [KW]
                                                        2.57 %
       EKWI
                       COILED COIL
                                                        2.33 %
```

SEQ MKDQQTVIMTECTSLQFVSPFAFEAMQKVDVVCLASLSDPELRLLLPCLVRMALCAPADQ

55

SEG

	WO	J 01/70434	PC1/1B01/02050
	MEM		••••••
	SEQ SEG	TSNLSANITGFLPIHCIYQLLRSRSFTKHKVSIKDWIYRQLCET	
5	PRD COILS	ccccceeeeeehhhhhhhhhhhhccccccchhhhhhhhh	ccccccccceeee
	MEM		
10	SEQ SEG	NSILTPASKSNPEATN@PVTE@EILNIF@GVIGGDNIRLN@RFS	• • • • • • • • • • • • • • • • • • • •
	PRD COILS	eeccccccccccchhhhhhhhhhcccccceeeehhh S	
15	MEM	••••••••••••	
	SEQ SEG	ALLANTKTLAAM@RKPKSYSSSLMD@IPIKFLIR@A@GL@@EL@	<xxxxxxxxx< p=""></xxxxxxxxx<>
20	PRD COILS	hhhhhhhhhhhhccccccccchhhhhhhhhhhhhhhhh	
	MEM	•••••••••••••••••••••••••••••••••••••••	
25	SEQ SEG	LCIVDDWICEEEITGTDALLRRMLLTNNAKNHSPKQLQEAFSAV	/PVNHTQVMQIIEHLTL
	PRD COILS	eeeecceeeechhhhhhhhhhhhhhcccccccchhhhhhh	
30	MEM		
	SEQ	LSASELIPYAEVLTSNMSQLLNSGVPRRILQTVNKLWMVLNTVM	
35	PRD COILS	hhhhhhhhhccccchhhhhccccchhhhhhhhhhhhhh	
<i></i>	MEM		
	SEG	KFVRQQKYTQNDLMIDPLIVLRCDQRVHRCPPLMDITLHMLNGY	
40	PRD COILS	hhhhhhhcccccccceeeeeccccccccccceeecccccc	: hhhhhhhhhhhhhhh
	MEM	•••••••••••••••••••••••••••••••••••••••	• • • • • • • • • • • • • • • • • • • •
45	SEG	QDRPSQNNTIGLVGQTDAPEVTREELKNALLAAQDSAAVQILLE	
	PRD COILS		
50	MEM	•••••••••••••••••••••••••••••••••••••••	
	SEQ SEG	SLLRNVQSVITTSAPNKGMEEGEDNLLCNLREVQCLICCLLHQM	1YIADPNIAKLVHFQGY
55	PRD	cceeeeeeeeccccccccchhhhhhhhhhhhhhhhhhh	
	115 (1		

	w	01/98454				-	PCT/IB01	/02050	
	SEQ	PCELLPLT	VAGIPSMHICL	DFIPELIA	QPELEKQIF	AIQLLSHL	CIQYALPI	<slsva< td=""><td>RLA</td></slsva<>	RLA
	SEG		• • • • • • • • • • •						
	PRD		eeeccceeee	ehhhhhhhh	hhhhhhhhh	hhhhhhhhl	nhhhhcc	chhhhh	hhh
5	COIL		• • • • • • • • • •						
5	MEM		• • • • • • • • • • • • • • • • • • • •						
							• • • • • • •	• • • • • •	• • •
	SEQ	VNVMGTLL	TVLTQAKRYAF	FMPTLPSL	VSFCRAFPP	LYEDIMSL	LIQIGQV	AVGZAD	TQT
	SEG								
10	PRD		հիհիհիհիհի	hcccccc	eeeccccc	chhhhhhhh	nhhhhccl	hhhhc	ccc
	COIL								
		• • • • • • •	• • • • • • • • •	• • • • • • •	• • • • • • • •	• • • • • • •			
	MEM		• • • • • • • • • • • • • • • • • • • •	• • • • • • • •	• • • • • • • •	• • • • • • • •	• • • • • • •	• • • • • •	• • •
15	SEQ	דדמודום	RLQQIKEKPSO	HEATENDE	CVVNCCDNT	CCMDDDUA	CUCTER		MMG
IJ	SEG		VERRIVERES						
	PRD		hhhhhhhccc						
	COIL					-cccccee	eeeeeem	1111111111111	eee
20	MEM			• • • • • • • •	• • • • • • • •				
	25.4								
	SEQ	VZGI				•			
	SEG								
25	PRD COIL:	eccc eccc							
43	MEM.		•						
	JIE11	• • • •							
						•			
				•					
30			Pro	site for	DKFZphte:	s3_20h12	. 2		
	10029	פכר	167->189	LEUCTNE	770050		PDOCUU	770	
	P2001		692->714	LEUCINE	_ZIPPER				
	. 5006	·	· u -		_~***		ייייייייייייייייייייייייייייייייייייייי	15.1	
35									
	(No F	fam data	available	for DKFZ	ohtes3_201	h]2.5)			

DKFZphtes3_21k14

5 group: testis derived

DKFZphtes3_21k14 encodes a novel 558 amino acid protein without similarity to known proteins.

10 No informative BLAST results; No predictive prosite, pfam or SCOP motife.

The new protein can find application in studying the expression profile of testis-specific genes.

15

unknown protein

perhaps complete cds.

20

Sequenced by LMU

Locus: unknown

25 Insert length: 2547 bp
Poly A stretch at pos. 2506, polyadenylation signal at pos. 2479

1501 TTCTGAATCA TCACTGGGAG CAAAACACAG ACTCACAGAG GAAGGGCAAG 1551 AGAAGGGTAA AGAACAAGAG AGACCACCTG AGGCAGTGAG CAAGTTTGCA BEDD AAGCGGAACA ATGAAGAAAC TGTAATGTCA GCTAGAGACA GGTACTTGGC 1651 CAGGCAGATG GCGCGGGTTA ATGCAAAGAC CTATATTGAG AAAGAAGATG 1701 ATTGATGGCT ACCCCAAGAG AAAGATTTAA GGAAGCACAG AAAACTGTAA 5 1751 TTCCTGGAAC CTGCTGCGTA AAACCATAAA GGAGTGTGTT ACCAGTAGTT 1801 TGGAGGGCAT TTTTAAATTT ATTTTCAAAA TTTTAAGTTA AAAGTCAGTC 1851 TTACAGCTTG GATGTTTGGA TGTGGATGTT TGGCTGAATT TATATATAGT 1901 GTGTACTCAT CAATACCACA TTCTTTGTTG TATTCAAGAA CCGTTAAGAG 1951 TGTGCTAATT CCCTGTAGGT ACATAATGAG GAAAATTTGC TCCACTACAA 10 2001 CCATTAAAAA ATAATTTTGG CCAGATACGG TAGCTCGTGC CTGTAATACC 2051 AACATTTTGG GAGGCCAAGG CAGAAGGATA TTGAGGCTAG GCATTCAAGA 2101 CCAGCCTAGG CAGGATAATA AGACCTTGTC TCTATTTAAA AAACAAAAAG 15 20

BLAST Results

25

No BLAST result

30

Medline entries

No Medline entry

35

Peptide information for frame 2

40 ORF from 29 bp to 1702 bp; peptide length: 558 Category: similarity to unknown protein Classification: Nucleic acid management

1 MAIPGRQYGL ILPKKTQQLH PVLQKPSVFG NDSDDDDETS VSESLQREAA

51 KKQAMKQTKL EIQKALAEDA TVYEYDSIYD EMQKKKEENN PKLLLGKDRK

101 PKYIHNLLKA VEIRKKEQEK RMEKKIQRER EMEKGEFDDK EAFVTSAYKK

151 KLQERAEEEE REKRAAALEA CLDVTKQKDL SGFYRHLLNQ AVGEEEVPKC

201 SFREARSGIK EEKSRGFSNE VSSKNRIPQE KCILQTDVKV EENPDADSDF

251 DAKSSADDEI EETRVNCRRE KVIETPENDF KHHRSQNHSR SPSEERGHST

50 301 RHHTKGSRTS RGHEKREDQH QQKQSRDQEN HYTDRDYRKE RDSHRHREAS

351 HRDSHWKRHE QEDKPRARDQ RERSDRVWKR EKDREKYSQR EQERDRQQND

401 QNRPSEKGEK EEKSKAKEEH MKVRKERYEN NDKYRDREKR EVGVQSSERN

451 QDRKESSPNS RAKDKFLDQE RSNKMRNMAK DKERNQEKPS NSESSLGAKH

501 RLTEEGQEKG KEQERPPEAV SKFAKRNNEE TVMSARDRYL ARQMARVNAK

10

BLASTP hits

No BLASTP hits available

5 Alert BLASTP hits for DKFZphtes3_21k14, frame 2

No Alert BLASTP hits found

Pedant information for DKFZphtes3_21k14, frame 2

Report for DKFZphtes3_21k14-2

15 **ELENGTHI** 567 67262.89 EMWI [pI] 8-96 TREMBL: ACOOL233_14 gene: "Fl2K2.14"; Arabidopsis []OMOL] thaliana chromosome II BAC F12K2 genomic sequence, complete sequence. 3e-11 20 **EFUNCATI** 04.99 other transcription activities ES. cerevisiae, YKR092cl le-05 **EFUNCATE** 30-10 nuclear organization ES. cerevisiae, YKR092c3 le-05 25 EFUNCATI 06.07 protein modification (glycolsylation, acylation, myristylation, palmitylation, farnesylation and processing) ES. cerevisiae, YKL201cl le-04 PF00748F **EBFOCKZ** BLO1182E Glycosyl hydrolases family 35 **EBFOCK21** 30 EEC] 2.7.1.37 Protein kinase 7e-06 [EC] 5.99.1.2 DNA topoisomerase 4e-Ob EPIRKW] phosphotransferase 7e-Ob pre-mRNA splicing le-Ob **EPIRKWI** citrulline 3e-0b **EPIRKUJ** tandem repeat 3e-06 35 **CPIRKWI** DNA binding 4e-06 **EPIRKUI** DNA replication 4e-Ob **EPIRKW CPIRKU** isomerase 4e-Db **EPIRKUD** ATP 3e-06 40 **EPIRKW3** phosphoprotein le-Ob -calcium binding 3e-06 [PIRKW] **EPIRKWI** alternative splicing 7e-06 **EPIRKUJ** P-loop 3e-06 **TPIRKUJ** EF hand 3e-Ob 45 **CPIRKUI** hair 3e-06 DEAD/H box helicase homology 3e-Ob **ESUPFAMI** unassigned Ser/Thr or Tyr-specific protein kinases 4e-**ESUPFAMI** calmodulin repeat homology 3e-Ob ESUPFAMI 50 CSUPFAMI unassigned ribonucleoprotein repeat-containing proteins le-Ob unassigned DEAD/H box helicases 3e-Ob **ESUPFAMD ESUPFAMI** trichohyalin 3e-06 **ESUPFAMI** protein kinase homology 4e-06 55 **ESUPFAM3** eukaryotic type I DNA topoisomerase 4e-06 ribonucleoprotein repeat homology le-Ob **ESUPFAMI** [KW] All_Alpha 22.75 % [KW] LOW_COMPLEXITY

	ZEG	### 7 PH P P P P P P P P P P P P P P P P P
5	PRD	ccccccccccccccccccccccccccccccccchhhhhh
•		
	SEQ	KQAMKQTKLEIQKALAEDATVYEYDSIYDEMQKKKEENNPKLLLGKDRKPKYIHNLLKAV
	SEG	······································
10	PRD	hhhhhhhhhhhhhhhcccccccchhhhhhhhhhhhhhhh
10	SEQ	EIRKKEQEKRMEKKIQREREMEKGEFDDKEAFVTSAYKKKLQERAEEEEREKRAAALEAC
	SEG	XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
	PRD	hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
15	SEQ	LDVTKQKDLSGFYRHLLNQAVGEEEVPKCSFREARSGIKEEKSRGFSNEVSSKNRIPQEK
15	SEG	FDALKKYDE201 LUHELINGVAGETETALKC21 KENKZGIKETKZVRI ZNEAZZVNIVILMEK
	PRD	hhhhhhhccchhhhhhhhhhhhhhhhhhhhhhhhhhhhh
20	SEG	CILQTDVKVEENPDADSDFDAKSSADDEIEETRVNCRREKVIETPENDFKHHRSQNHSRS
20	PRD	hhhhhhhhhhccchhhhhhhhhhhhhhhhhhhhhhhhh
•	FILD	THE TAX TO
	SEQ	PSEERGHSTRHHTKGSRTSRGHEKREDQHQQKQSRDQENHYTDRDYRKERDSHRHREASH
	SEG	***************************************
25	PRD	cccchhhhhhhhhhcccchhhhhhhhhhhhhhhhhhhhh
	SEQ	RDSHWKRHEQEDKPRARDQRERSDRVWKREKDREKYSQREQERDRQQNDQNRPSEKGEKE
	SEG	xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
20	PRD	hhhhhhhhccccchhhhhhhhhhhhhhhhhhhhhhhhhh
30	SEQ	EKSKAKEEHMKVRKERYENNDKYRDREKREVGVQSSERNQDRKESSPNSRAKDKFLDQER
	SEG	XXXXXXXXX
	PRD	hhhhhhhhhhhhhhhccchhhhhhhhhhhhhhhhhhhhh
25	SE4	CHICHDING A CARCENIA CARCACTORIA CARCACTOR
35	SEQ	SNKMRNMAKDKERNQEKPSNSESSLGAKHRLTEEGQEKGKEQERPPEAVSKFAKRNNEET
	PRD	hhhhhhhhhhhhhhccchhhhhhhhhhhhhhhhhhccccc
40	SEG	VMSARDRYLARQMARVNAKTYIEKEDD
40	SEG PRD	hhhhhhhhhhhhhhhhhhccc
	FILD	THE THE PARTY OF T
	4 NI	D. Charles and Della Con DEET Land Della Della
45	ONI	Prosite data available for DKFZphtes3_21k14.2)
	(No	Pfam data available for DKFZphtes3_21k14.2)

PCT/IB01/02050

DKFZphtes3_22ill

WO 01/98454

10

20

25

30

5 group: testis derived

DKFZphtes3_22ill encodes a novel 580 amino acid protein with similarity to RCCl-like G exchanging factor RLG, UVRA (UVB-resistance protein) of Arabidopsis thaliana and to the murine retinitis pigmentosa GTPase regulator.

No informative BLAST results; No predictive prosite, pfam or SCOP motife.

15 The new protein can find application in studying the expression profile of testis-specific genes.

Homo sapiens chromosome 7q22 sequence, ORF4, extension

differences to genmodel of ORF4, differential splicing

Sequenced by LMU

Locus: /map="?q22"

Insert length: 2236 bp

Poly A stretch at pos. 2197, polyadenylation signal at pos. 2180

L ACANTGCTCA GATCGGGAGG TGGAGCCAAT CAGGTCCAAC CAAGAGGAGG 51 GGACACCGGC ACTCCACTAG CAGGAAAACG GGCCGAGGGA CCGCAAGCAG 101 GGGGTGCCTA GTCCTCGTCC CCCAAAGACC AATCGTAAGC CAGATACAGG 35 151 CGAGTGACTG TCAAGAAGGC CAATTAGAGC CTCCGAAGGG AATCTGGACC 201 TGCCTCTTCT CTGAGGGACG GCTCTACCTA CCAATAGCAT GGGCGAGAAG 251 GCGGTCCCTT TGCTAAGGAG GAGGCGGGTG AAGAGAAGCT GCCCTTCTTG **JUL TGGCTCGGAG CTTGGGGTTG AAGAGAGAG GGGGAAAGGA AATCCGATTT** 351 CCATCCAGTT GTTCCCCCCA GAGCTGGTGG AGCATATCAT CTCATTCCTC 40 4D1 CCAGTCAGAG ACCTTGTTGC CCTCGGCCAG ACCTGCCGCT ACTTCCACGA 451 AGTGTGCGAT GGGGAAGGCG TGTGGAGACG CATCTGTCGC AGACTCAGTC 5D1 CGCGCCTCCA AGATCAGGGT TCTGGAGTCC GGCCCTGGAA GAGAGCTGCC 551 ATTCTGAACT ACACGAAGGG CCTGTATTTC CAGGCATTTG GAGGCCGCCG LOD CCGATGTCTC AGCAAGAGCG TGGCCCCCTT GCTAGCCCAC GGCTACCGCC 651 GCTTCTTGCC CACCAAGGAT CACGTCTTCA TTCTTGACTA CGTGGGGACC 45 7D1 CTCTTCTTCC TCAAAAATGC CCTGGTCTCC ACCCTCGGCC AGATGCAGTG 751 GAAGCGGGCC TGTCGCTATG TTGTGTTGTG TCGTGGAGCC AAGGATTTTG BD1 CCTCGGACCC AAGGTGTGAC ACAGTTTACC GTAAATACCT CTACGTCTTG 851 GCCACTCGGG AGCCGCAGGA AGTGGTGGGT ACCACCAGCA GCCGGGCCTG 50 9D1 TGACTGTGTT GAGGTCTATC TGCAGTCTAG TGGGCAGCGG GTCTTCAAGA 951 TGACATTCCA CCACTCAATG ACCTTCAAGC AGATCGTGCT GGTTGGTCAG LDDL GAGACCCAGC GGGCTCTACT GCTCCTCACA GAGGAAGGAA AGATCTACTC 1051 TTTGGTAGTG AATGAGACCC AGCTTGACCA GCCACGCTCC TACACGGTTC LIDL AGCTGGCCCT GAGGAAGGTG TCCCACTACC TGCCTCACCT GCGCGTGGCC
LISL TGCATGACTT CCAACCAGAG CAGCACCCTC TACGTCACAG ACCAGGGGGG
LZDL AGTGTATTTT GAGGTGCATA CCCCAGGGGT GTATCGCGAT CTCTTTGGGA
LZSL CCCTTCAAGC CTTTGACCCC CTGGACCAGC AGATGCCGCT TGCTCTCTCA
LJDL CTGCCTGCCA AGATCCTATT CTGTGCTCTT GGCTACAACC ACCTTGGCCT 55

	WO 01/98454				PCT/IB01/02050
	14D1 AGCTAGGAAC	TTTGGCCGAA AGGGGACAAA	ATGGACCGAG	AGGAAATAAC GGGAACCCAC	ACAGGTTTGT
_	1501 GCTGAGCCAG		TCAGCAAGGA	GCTGCTGGGC	TGCGGCTGTG
5	1551 GGGCTGGGGG 1601 AAGCTCCAAG 1651 GGAGTGCCTA	TCAAGGTCCC TACATCCTGT	TCTGTGTGCC CCAGCCACGA	CATTGAGCAG	GTGCCACCAG CACGCCCCCT
10	17D1 ATCGCCACCT 1751 GGGGCCAGAG 18D1 GTACCTCAGC	CACCCCAGCA CACCCCAGGA CAGATCCACA	CCCCGGGGGG	ATGGCCCAGG	CCTGCGAGGA
10	1851 AGATGAAGGA 1901 TTCTTCTGGG	GATCGTAGGG	TGGATGCCCC CATGCTGCAG	GTTGCAGGAC TGATGGCCGC AGGGCTGAAG	ACAGAAGGAC GAGGCGGGGG
15	2001 GTCCCTGGAG 2051 CCATTGTGCA 2101 AACCAGGGTA	GAGGGAGTCC CATGCGTGTG	GGCCCCAGGC GGAAGGGGTT	CAGGGACTAA GCTAGGGGGT	GGAGCAATGA GGGGACGGCT
	2151 ACTATCATGG	ACAAGAGATT	TGATGGATAG	AATAAAAGGC	TGCAGCGAAA
20			DIAST D	.2.5	
			BLAST Resu	11 C S	
25	Entry AF053356 for the Homo sapiens chroscore = 2952, P = 10 exons	omosome 7q2	sequence,	complete se = 666/729	equence.
30	•	-	Medline ent	ries	
35	No Medline entry				
					•
		Peptide	information	for frame	. -
40	ORF from 239 bp t Category: similar Classification: r	ity to unkr	peptide len nown protein	igth: 580	
45	1 MGEKAVPLLR 51 ISFLPVRDLV 101 KRAAILNYTK 151 YVGTLFFLKN	ALGQTCRYFH GLYFQAFGGR	EVCDGEGVWR RRCLSKSVAP	RGKGNPISIQ RICRRLSPRL LLAHGYRRFL	QDQGSGVRPW PTKDHVFILD
50	201 LYVLATREPQ 251 LVGQETQRAL	EVVGTTSSRA LLLTEEGKIY	CDCVEVYLQS SLVVNETQLD	SGQRVFKMTF QPRSYTVQLA	HHSMTFKQIV LRKVSHYLPH
	301 LRVACMTSNQ 351 LALSLPAKIL	SSTLYVTDQG	GVYFEVHTPG	VYRDLFGTLQ	AFDPLDGGMP

401 TQVCYLQRPI TLWCGLNHSL VLSQSSEFSK ELLGCGCGAG GRLPGWPKGS 451 ASFVKLQVKV PLCACALCAT RECLYILSSH DIEQHAPYRH LPASRVVGTP

501 EPSLGARAPQ DPGGMAQACE EYLSQIHSCQ TLQDRTEKMK EIVGWMPLMA

551 AQKDFFWEAL DMLQRAEGGG GGVGPPAPET

55

BLASTP hits

No BLASTP hits available

5

Alert BLASTP hits for DKFZphtes3_22ill, frame 2

TREMBL:AF053356_11 product: "ORF4": Homo sapiens chromosome 7q22 sequence: complete sequence: N = 1: Score = 1554: P = 1:6e-159

TREMBL:AF130441_1 gene: "UVRA"; product: "UVB-resistance protein UVRA";
Arabidopsis thaliana UVB-resistance protein UVRA (UVRA) mRNA;

Arabidopsis thaliana UVB-resistance protein UVRA (UVRA) mRNA, complete

15 cds., N = 1, Score = 109, P = 0.0082

Length = 318

TREMBL:AFO44677_1 gene: "Rpgr"; product: "retinitis pigmentosa GTPase

regulator"; Mus musculus retinitis pigmentosa GTPase regulator 20 (Rpgr)

mRNA, complete cds., N = 1, Score = 106, P = 0.035

>TREMBL:AF05335b_ll product: "ORF4"; Homo sapiens chromosome
25 7q22 sequence; complete sequence.

HSPs:

30

40

Score = 1554 (233.2 bits), Expect = 1.6e-159, P = 1.6e-159
Identities = 303/318 (95%), Positives = 303/318 (95%)

Query: 1
35 MGEKAVPLLRRRRVKRSCPSCGSELGVEEKRGKGNPISIQLFPPELVEHIISFLPVRDLV 60

MGEKAVPLLRRRRVKRSCPSCGSELGVEEKRGKGNPISIQLFPPELVEHIISFLPVRDLV Sbjct: 1
MGEKAVPLLRRRRVKRSCPSCGSELGVEEKRGKGNPISIQLFPPELVEHIISFLPVRDLV 60

Query: 61
ALGQTCRYFHEVCDGEGVWRRICRRLSPRLQDQGSGVRPWKRAAILNYTKGLYFQAFGGR 120
ALGQTCRYFHEVCDGEGVWRRICRRLSPRLQDQ
TKGLYFQAFGGR

45 Sbjct: 61 ALGQTCRYFHEVCDGEGVWRRICRRLSPRLQDQD------TKGLYFQAFGGR 106

Query: 121 RRCLSKSVAPLLAHGYRRFLPTKDHVFILDYVGTLFFLKNALVSTLGQMQWKRACRYVVL 180

50
RRCLSKSVAPLLAHGYRRFLPTKDHVFILDYVGTLFFLKNALVSTLGQMQWKRACRYVVL
Sbjct: 107
RRCLSKSVAPLLAHGYRRFLPTKDHVFILDYVGTLFFLKNALVSTLGQMQWKRACRYVVL 166

55 Query: 181
CRGAKDFASDPRCDTVYRKYLYVLATREPQEVVGTTSSRACDCVEVYLQSSGQRVFKMTF 240
CRGAKDFASDPRCDTVYRKYLYVLATREPQEVVGTTSSRACDCVEVYLQSSGQRVFKMTF

Sbjct: 167

CRĞAKDFASDPRCDTVYRKYLYVLATREP@EVVGTTSSRACDCVEVYL@SSG@RVFKMTF 226

Query: 241

5 HHSMTFKQIVLVGQETQRALLLLTEEGKIYSLVVNETQLDQPRSYTVQLALRKVSHYLPH 300

HHSMTFKQIVLVGQETQRALLLLTEEGKIYSLVVNETQLDQPRSYTVQLALRKVSHYLPH

Sbjct: 227

HHSMTFKQIVLVGQETQRALLLLTEEGKIYSLVVNETQLDQPRSYTVQLALRKVSHYLPH 286

10

Query: 301 LRVACMTSNQSSTLYVTD 318

LRVACMTSN@SSTLYVTD

Sbjct: 287 LRVACMTSNQSSTLYVTD 304

15

Pedant information for DKFZphtes3_22ill, frame 2

Report for DKFZphtes3_22ill-2

20

ELENGTHD 580

EMW3 64889-49

[pI] 9.01

25 EHOMOLI TREMBL:AFD5335b_11 product: "ORF4": Homo sapiens chromosome 7q22 sequence. complete sequence. le-174
EBLOCKSI BLOOb25B Regulator of chromosome condensation (RCC1)

proteins

EBLOCKSI BLOOBESA Regulator of chromosome condensation (RCCL)

30 proteins

SEQ SEG PRD

[KW] Alpha_Beta

EKWI LOW_COMPLEXITY 3.62 %

35 MGEKAVPLLRRRRVKRSCPSCGSELGVEEKRGKGNPISIQLFPPELVEHIISFLPVRDLV ZEQ SEG PRD SEQ ALGQTCRYFHEVCDGEGVWRRICRRLSPRLQDQGSGVRPWKRAAILNYTKGLYFQAFGGR 40 SEG PRD SEQ RRCLSKSVAPLLAHGYRRFLPTKDHVFILDYVGTLFFLKNALVSTLG@M@WKRACRYVVL SEG 45 PRD SEQ CRGAKDFASDPRCDTVYRKYLYVLATREP@EVVGTTSSRACDCVEVYL@SSG@RVFKMTF SEG PRD 50 HHSMTFKQIVLVGQETQRALLLLTEEGKIYSLVVNETQLDQPRSYTVQLALRKVSHYLPH SEQ SEG PRD 55 LRVACMTSNQSSTLYVTDQGGVYFEVHTPGVYRDLFGTLQAFDPLDQQMPLALSLPAKIL

	V	WO 01/98454	PCT/IB01/02050
	SEQ	The state of the s	
	PRD	· · · · · · · · · · · · · · · · · · ·	
5	SEQ	The state of the s	
	SE <i>G</i> PRD		
	SEQ	The same of the sa	
10	SEG PRD		
	SEQ	₹ EIVGWMPLMAA@KDFFWEALDML@RAEGGGGGVGPPAPET	
15	SEG PRD		
	(No	o Prosite data available for DKFZphtes3_22ill.2)
20	(No	Pfam data available for DKFZphtes3_22ill-2)	

5 group: testis derived

DKFZphtes3_22124 encodes a novel 451 amino acid protein with similarity to the F-box protein FBL2 of the rat.

No informative BLAST results; No predictive prosite; pfam or SCOP motife.

The new protein can find application in studying the expression profile of testis-specific genes.

15

similarity to p37NB (Homo sapiens)

Sequenced by LMU

20

Locus: /map="7q22-q31.1"

Insert length: 1537 bp

Poly A stretch at pos- 1459, no polyadenylation signal found

25

	7	CAACAGGACG	ATGCGACTCC	TGCCGAGGCA	CTTCCACAAC	TTACAGAATC
	51	TTAGTTTGGC		CGGTTCACAG	ACAAAGGCTT	ACAGTACCTG
	101	AACTTGGGGA	ATGGATGCCA	CAAGCTCATC	TATCTGGACC	TCTCTGGCTG
30	151	CACCCAGATT	TCAGTCCAAG	GCTTCAGGTA	CATTGCAAAC	AGCTGCACTG
	507	GAATTATGCA	TCTTACCATT	AATGACATGC	CAACTCTGAC	GGACAACTGT
	527	GTAAAAGCTT	TAGTTGAAAA	ATGCTCTCGT	ATTACATCGC	TGGTTTTCAC
	307	TGGTGCACCG	CATATCTCCG	ATTGTACTTT	CAGAGCTCTT	TCTGCTTGTA
	351	AACTCAGAAA	GATCCGATTT	GAAGGAAATA	AAAGGGTTAC	TGATGCATCC
35	401	TTCAAATTTA	TAGACAAGAA	TTATCCAAAT	CTCAGTCACA	TTTATATGGC
	451	TGACTGCAAG	GGAATAACAG	ACAGCAGCCT	CAGATCCCTT	TCACCTTTGA
	5D1	AGCAACTGAC	TGTGTTGAAT	TTGGCAAATT	GTGTAAGAAT	TGGTGATATG
	551	GGACTAAAGC	AATTTCTTGA	TGGTCCTGCA	AGCATGAGGA	TAAGAGAGCT
	POT	AAATTTAAGC	AACTGTGTGC	GGCTAAGTGA	TGCCTTTGTT	ATGAAACTAT
40	65]	CTGAGCGCTG	CCCTAATTTA	AACTACTTGA	GTTTACGAAA	TTGTGAACAT
	701	TTGACTGCCC	AAGGAATTGG	ATATATTGTA	AACATCTTTT	CCTTGGTATC
	751	AATAGATCTC	TCTGGAACAG	ACATCTCTAA	TGAGGGTTTG	AATGTGCTTT
	807	CCAGACATAA	AAAATTGAAG	GAACTTTCTG	TATCTGAATG	TTATAGAATC
	851	ACTGATGATG	GAATTCAGGC	ATTCTGCAAA	AGCTCACTGA	TCTTGGAACA
45	401	TTTGGATGTC	TCTTATTGCT	CCCAGCTGTC	AGATATGATT	ATCAAAGCAC
	951	TGGCCATTTA	CTGCATTAAC	CTCACATCTC	TCAGCATTGC	TGGCTGTCCA
	7007	AAGATTACTG	ACTCAGCAAT	GGAGATGTTA	TCGGCAAAAT	GCCATTACCT
	1051	GCACATTTTG	GATATCTCTG	GTTGTGTCTT	GCTTACTGAC	CAAATCCTTG
	7707	AGGACCTTCA	GATAGGCTGC	AAACAACTCC	GGATCCTTAA	GATGCAATAC
50	1151	TGCACAAATA	TTTCCAAGAA	GGCAGCTCAA	AGAATGTCAT	CTAAAGTTCA
	7507	GCAGCAGGAA	TACAACACTA	ATGACCCTCC	ACGTTGGTTT	GGCTATGATA
	1251	GGGAAGGAAA	CCCTGTTACA	GAGCTTGACA	ACATAACATC	ATCTAAAGGA
	7307	GCCTTAGAAT	TAACAGTGAA	AAAGTCAACA	TACAGCAGTG	AAGACCAAGC
	1351	AGCGTGACCT	TCAGCCTCAA	GCAGGAAGAA	CAAAAAATCA	AGAACTTGGC
55	1401	AAGTTTTCTC	CATTTGTTGC	AAGTATGTTT	ACTAGCTGAA	TCTCAATAAC
	1451	AATGTAAACA	AGCAAAAAA	AAAAAAAAA	AAAAAAAAA	AAAAAAAAA
	1501	AAAAAAAAA		AAAAAAAAA	AAAAAG	

BLAST Results

- 5 Entry ACOO5250 from database EMBL:
 Homo sapiens BAC clone R6318MO5 from 7q22-q31.1, complete
 sequence.
 Score = 830, P = 1.8e-124, identities = 180/193
- 10 Entry HS32907 from database EMBL:
 Human p37NB mRNA, complete cds.
 Score = 318, P = 4.6e-04, identities = 70/78

15

Medline entries

97136875:

- 20 Kim D. LaQuaglia MP. Yang SY.; A cDNA encoding a putative 37 kDa leucine-rich repeat (LRR) protein. p37NB. isolated from S-type neuroblastoma cell has a differential tissue distribution. Biochim Biophys Acta 1996
- 25 Dec 11:1309(3):183-8

Peptide information for frame 2

ORF from 11 bp to 1354 bp; peptide length: 448 Category: similarity to known protein Classification: unclassified

1 MRLLPRHFHN LQNLSLAYCR RFTDKGLQYL NLGNGCHKLI YLDLSGCTQI
51 SVQGFRYIAN SCTGIMHLTI NDMPTLTDNC VKALVEKCSR ITSLVFTGAP
101 HISDCTFRAL SACKLRKIRF EGNKRVTDAS FKFIDKNYPN LSHIYMADCK
40 151 GITDSSLRSL SPLKQLTVLN LANCVRIGDM GLKQFLDGPA SMRIRELNLS
201 NCVRLSDAFV MKLSERCPNL NYLSLRNCEH LTAQGIGYIV NIFSLVSIDL
251 SGTDISNEGL NVLSRHKKLK ELSVSECYRI TDDGIQAFCK SSLILEHLDV
301 SYCSQLSDMI IKALAIYCIN LTSLSIAGCP KITDSAMEML SAKCHYLHIL
351 DISGCVLLTD QILEDLQIGC KQLRILKMQY CTNISKKAAQ RMSSKVQQQE
45 401 YNTNDPPRWF GYDREGNPVT ELDNITSSKG ALELTVKKST YSSEDQAA

BLASTP hits

50

35

No BLASTP hits available

Alert BLASTP hits for DKFZphtes3_22124, frame 2

55 No Alert BLASTP hits found

Pedant information for DKFZphtes3_22124, frame 2

Report for DKFZphtes3_22124.2

```
ELENGTHD 451
   EMUD 50545.95
Epid 8.68
   EHOMOLI
               TREMBLNEW: AF186273_1 product: "leucine-rich
   repeats containing F-box protein FBL3"; Homo sapiens leucine-rich repeats containing F-box protein FBL3 mRNA; complete cds. &e-31
   YJR090c3 8e-20
   EFUNCATI 03.04 budding, cell polarity and filament formation ES. cerevisiae, YJR090cl 8e-20
EFUNCATI 01.05.04 regulation of carbohydrate utilization E
   cerevisiae, YJR090cl 8e-20
20
   EFUNCATI 11.04 dna repair (direct repair, base excision repair
   EFUNCATI 30-10 nuclear organization ES- cerevisiae, YJRD52wl
   3e-07
   BLOCKZI PRODOJAB
BBLOCKZI PRODOJAD
25
   EBLOCKSI ROUSELE
EBLOCKSI BPD1921A
   EPIRKW3
               tandem repeat 2e-18
   EPIRKWI
EPIRKWI
               zinc finger le-07
DNA binding le-07
30
   ISUPFAMI leucine-rich alpha-2-glycoprotein repeat homology 2e-18.
   [SUPFAM] regulatory protein ESAG&c le-07
           Alpha_Beta
   EKW3
35
   SEQ
       NRTMRLLPRHFHNLQNLSLAYCRRFTDKGLQYLNLGNGCHKLIYLDLSGCTQISVQGFRY
   PRD
       IANSCTGIMHLTINDMPTLTDNCVKALVEKCSRITSLVFTGAPHISDCTFRALSACKLRK
40
   SEQ
   PRD
      SEQ IRFEGNKRVTDASFKFIDKNYPNLSHIYMADCKGITDSSLRSLSPLKQLTVLNLANCVRI
   PRD
       45
       GDMGLKQFLDGPASMRIRELNLSNCVRLSDAFVMKLSERCPNLNYLSLRNCEHLTAQGIG
   SEQ
   PRD
       cccccccccccccccchhhhhhccccccccccccccccee
   SEQ
       YIVNIFSLVSIDLSGTDISNEGLNVLSRHKKLKELSVSECYRITDDGIQAFCKSSLILEH
50
       PRD
       LDVSYCSQLSDMIIKALAIYCINLTSLSIAGCPKITDSAMEMLSAKCHYLHILDISGCVL
   SEQ
       PRD
      LTDQILEDLQIGCKQLRILKMQYCTNISKKAAQRMSSKVQQQEYNTNDPPRWFGYDREGN
55
   SEQ
      PRD
   SEQ PVTELDNITSSKGALELTVKKSTYSSEDQAA
```

PRD ccccccccceeeeeccccccccc

5

(No Prosite data available for DKFZphtes3_22124.2)

(No Pfam data available for DKFZphtes3_22124.2)

DKFZphtes3_26g3

5 group: testis derived

DKFZphtes3_26g3 encodes a novel 1090 amino acid protein without similarity to known proteins.

No informative BLAST results; No predictive prosite, pfam or SCOP motife.

The new protein can find application in studying the expression profile of testis-specific genes.

15

similarity to C.elegans CO9D4.4

on genomic level encoded by HSDJ19819 20 perhaps complete cds.

Sequenced by EMBL

Locus: /map="b"

25

Insert length: 4562 bp Poly A stretch at pos. 4550, polyadenylation signal at pos. 4565

30 1 GATTCAGTTA CTGAAGACTT AGATGCACCC TGGATGGGAA TTCAGAATCT 51 TCAGAGATCA GAGTCCAGTA AAATGGATAA ATATGAGACT GAAGAAAGCT LOL CTGTAGCAGG ACTTTCTAGC CCAGAGTTGA AAGTCAGACC TGCTGGTGCC
LSL TCCAGTATTT GGTATACAGA AGGTGAAAAG CAGCTAACAA AATCTCTAAA
LOL AGGAAAGAAT GAAGAATCAA ATAAATCCAA AGTTAAGGTT ACTAAGCTTA
LSL TGAAAACAAT GAAATCTGAA AACACAAAAA AATTAATAAA ACAGAACTCT 35 BOL AAGGATTCTG TGGTTTTGGT AGGCTACAAA TGTTTGAAAA GTACAGCATC 351 AAATGATCTC ATTAAATGCT TTGAAGGCAA TCCTTCACAT AGTCAGAAGG 451 ATGAGACAGA CAAGTCAAAA GGAAGCTAGC TGTTTGCCAA CTAATACAGA 40 501 GAGAACTGAA CAAAAGTCTC CAGATATTGA AAATGTTCAA CCAGACCAGT 551 TTGATCCTTT GAACTCTGGC AACCTAAATC TTTGTGCAAA TTTGTCCATT 601 TCAGGTAAAC TTGATATCTC CCAGGACGAT AGTGAAATTA CACAAATGGA LSI ACACAATCTG GCATCCAGAA GGTCATCAGA CGATTGCCAT GATCATCAAA 701 CAACCCCATC TTTGGGAGTT AGAACAATTG AAATAAAGCC CAGTAATAAA
751 GATCCTTTCA GTGGAGAGAA TATAACTGTC AAACTAGGAC CTTGGACAGA
801 GCTTCGACAA GAGGAAATAC TTGTGGATAA TTTACTACCC AACTTTGAGT
851 CCTTAGAATC TAATGGTAAA TCTAAATCTA TAGAAATAAC ATTTGAAAAG 45 901 GAAGCTTTGC AAGAAGCAAA GTGTCTTTCT ATTGGAGAAT CATTAACTAA
951 ATTACGAAGT AATCTACCTG CCCCTTCTAC AAAAGAATAT CATGTTGTAG
1001 TAAGTGGAGA TACAATTAAG TTACCAGATA TTAGTGCCAC ATATGCCTCA
1051 TCTAGATTTT CAGATTCAGG TGTTGAAAGT GAACCGAGTT CTTTTGCGAC
1101 ACATCCAAAC ACTGATTTAG TCTTTGAAAC TGTGCAAGGG CAAGGTCCTT
1151 GCAATAGTGA AAGATTATTT CCTCAGCTTT TGATGAAACC TGATTATAAT 50 1201 GTAAAATTTT CATTAGGAAA TCATTGTACT GAGAGTACAA GTGCTATAAG 1251 TGAAATACAG TCATCTTTGA CATCCATAAA CTCTCTACCC TCCGATGATG
1301 AACTGTCACC TGATGAAAAT TCTAAGAAAT CTGTTGTACC TGAATGCCAT
1351 CTAAATGATA GCAAAACTGT ATTAAATCTA GGAACGACTG ATTTGCCAAA 55 14D1 ATGTGATGAT ACTAAAAAGT CAAGTATCAC TTTGCAACAG CAGAGTGTTG

						_
	1451	TATTTTCAGG	GAACTTGGAC	AATGAAACTG	TAGCAATACA	TTCCTTAAAT
	1501	TCAAGCATTA	AAGACCCTTT	ACAATTTGTT	TTTTCAGATG	AAGAGACTTC
	1551	CAGTGATGTG	AAAAGTAGTT	GCAGCTCCAA	ACCTAACTTG	GATACTATGT
	7207	GTAAAGGCTT	CCAGAGTCCT	GATAAATCTA	ATAACTCTAC	AGGGACAGCA
5	1651	ATTACATTAA	ATTCAAAACT	GATTTGTTTA	GGCACTCCTT	
5						GTGTCATTTC
	1701	AGGTTCCATT	TCTAGTAATA	CAGATGTTAG	TGAAGATAGA	ACTATGAAAA
	1751	AAAATAGTGA	TGTATTAAAT	CTCACACAGA	TGTATTCAGA	AATCCCTACA
	reor	GTTGAAAGTG	AAACTCATCT	GGGTACAAGT	GATCCTTTTT	CAGCCAGTAC
	1851	TGATATAGTA	AAGCAAGGGC	TTGTGGAAAA	TTATTTTGGT	TCTCAAAGCA
10	1901	GTACGGATAT	TTCTGACACA	TGTGCTGTTA	GCTACAGCAA	TGCACTTAGC
	1951	CCTCAGAAGG	AAACTTCTGA	AAAAGAAATT	AGTAATCTTC	AGCAGGAACA
	2007	GGATAAAGAG	GATGAGGAGG	AAGAGCAGGA	TCAACAAATG	GTTCAAAATG
				· · · · · · · · · · · · · · · · · · ·		
	2051	GGTACTATGA	AGAAACAGAT	TATTCAGCTT	TGGATGGAAC	AATAAATGCT
	5707	CACTATACAA	GCAGAGATGA	ACTAATGGAA	GAAAGACTTA	CAAAATCTGA
15	2151	AAAAATAAAC	AGTGACTATC	TGAGAGATGG	TATAAACATG	CCTACTGTCT.
	5507	GTACTTCTGG	TTGTTTGTCC	TTCCCGTCTG	CACCACGAGA	GTCTCCTTGT
	2251	AATGTTAAAT	ATTCTTCCAA	AAGTAAATTT	GATGCCATTA	CAAAGCAGCC
	5301	AAGCAGTACT	TCTTACAACT	TCACTTCTTC	GATTTCCTGG	TATGAAAGTT
	2351	CACCAAAACC	TCAAATACAA	GCCTTCCTTC	AGGCAAAAGA	AGAACTGAAG
20	2401	CTACTAAAAC	TTCCTGGGTT	CATGTACAGT	GAAGTTCCTC	TGCTGGCATC
20						
	2451	CTCAGTACCT	TATTTTAGTG	TAGAAGAAGA	GGGTGGTTCT	GAAGATGGAG
	2501	TACATCTGAT	TGTCTGTGTG	CACGGTTTAG	ATGGAAACAG	TGCAGATCTC
	2551	CGATTAGTAA	AAACTTACAT	TGAACTTGGA	TTGCCTGGGG	GAAGAATTGA
	5P07	TTTTCTTATG	TCTGAGAGAA	ATCAGAATGA	TACTTTTGCT	GATTTTGATA
25	2651	GCATGACTGA	TCGTCTTTTG	GATGAGATAA	TACAGTATAT	TCAGATATAT
	2701	AGTCTAACAG	TCTCAAAAAT	AAGCTTTATT	GGACATTCGT	TGGGCAATTT
	2751	AATAATTCGT	TCAGTGCTTA	CAAGGCCAAG	GTTTAAATAT	TACCTCAACA
	2801	AACTTCATAC	CTTTCTGTCT	CTTTCTGGAC	CTCACCTTGG	TACACTCTAC
	2851	AACAGCAGTG	CTCTTGTTAA	TACAGGTCTC	TGGTTTATGC	AGAAATGGAA
20		AAAATCAGGT	TCGCTTTTGC	AGCTGACATG	TCGAGATCAC	
30	2901					TCAGACCCTC
	2951	GCCAAACTTT	TTTATATAAG	CTTAGTAACA	AAGCAGGGCT	TCATTATTTC
	3007	AAAAATGTTG	TGCTAGTGGG	ATCCCTACAG	GATCGCTATG	TTCCTTATCA
	3051	CTCTGCCCGC	ATTGAAATGT	GTAAAACAGC	TTTAAAGGAC	AAACAGTCAG
	3707	GACAGATCTA	TTCAGAAATG	ATCCACAACT	TGCTTCGACC	CGTTCTGCAA
35 ·	3151	AGCAAGGACT	GTAATTTGGT	TCGCTATAAT	GTCATCAATG	CATTGCCCAA
	3507	TACAGCTGAT	TCACTCATTG	GGAGAGCTGC	ACATATAGCT	GTTCTTGATT
	3251	CGGAAATATT	TTTAGAGAAA	TTCTTTCTGG	TTGCTGCCCT	CAAATATTTC
	3307	CAATAGTATA	AAAGCATTGT	TAGCGACTGG	ACAATTACCT	CATTCAACAA
	3351	TGTTTCAAAT	AATGTATTAT	ATTAAAATGT	AGATGCTGAT	AAGTTCTAAG
40		AAATATTTAT	ACCTTTTTAT	ATGGAAGATA	ATTTATATCA	_
40	3401					TCCATGTTTA
• • •	3451	GTGCTTTTTA	AACATCAACT	TTACTTTCTA	GGTAATGTGG	CTGTGCAATA
	3501	TTTTTTTAAT	TTTATCTTTT	TACTTTTCTA	TTACTTTTTC	ATATATTTTG
	3551	CTACCTAAGT			CCATACCTGT	GTCTGATTGT
	3203	TTATTATTGG	CTTTCCACAA	TTCTTACATC	AGACTACATT	ATATTAGAGA
45	3651	CCATTATTGC	TAGAATAGCA	TGGGATTTAA	AATTTCTAAT	ACTGGGGGTA
	3701	TTATTTAGTT	TAAATTAAAT	TTTTCTTTTC	ACATTTTACT	GTGTTTTAAC
	3751		ATTATGGCTG	CTACAATATA	TTTTTTGAAA	TCAACTTCTG
	3801	TAGTTCTAAA	ATACAACTTT	ATCATACAAT	CAAACCAGGT	AGTTCATATA
	3851	AAACAGTGTA	ATACAAGTTT	TCTATAAAGT	CATTACTETT	GCTTAAACAT
50	3901	ATTTCATGCC	TATTAAAATA	TATTTTCTAC	TGGTGATTTC	AACATTATTT
	3951	CTCATACTGA	CTTTTATTAC	TGGAAATGTT	CCTGTACATG	TTGGCAGCAG
	4001	ATAAAGATTT	TTGAATGTTT	GAATGCCCTC	TGCCTTGATT	TGGTTGGATT
	4051	TTGCTAATTG	GTATGTTGCT	TGAACTTTAT	GACTACATTT	TCTTTTAACT
	4101	TTTTTCATGG	ACTTCCTTAT	ATGTACATAA	TAATTAAATG	TTGAAATTTA
55	4151	TGAAATACTT	TTATGAATTT	AGATAATTTT	TAAATATTGT	TAAAATTTAT
<i>JJ</i>	4507	TGAACTAAAA	AGTAATGTAA	TAATAAATA	TCATGTTAAA	GATGGAACAA
	4251	AATAATTAAC	TTTACATGTT	TGGTGATACA	GATGCAAATG	TTTTTGATAT
	4301	ATGGAGATGT	TGAGTCTTTT	GACTTTACTA	AAGGTGCTGA	ATAGCATTAA

4351 ATTCACTATT TTCCTTTTCT GTTTTACTTG TGAAAATAAA AATGCACTAA
4401 GGTTGGGTAG AAGTTCTGTT TGCACTCACT AATTGTGACA GACAGAGGTT
4451 TTTGTAAGTA TTTATTGTAC AATTGATGCA TGTTTATTTT TAGCGTTGTT
4501 ATTGCCTCTG GTGTTAATAA ATGAACAAAT GGCTATCTGG AGGAACAGCT
5 4551 AAAAAAAAA AA

BLAST Results

10

Entry HSDJ19819 from database EMBLNEW:
Human DNA sequence *** SEQUENCING IN PROGRESS *** from clone
DJ19819
Score = 7221, P = 0.0e+00, identities = 1455/1461

15

Medline entries

20

No Medline entry

25

30

Peptide information for frame 1

ORF from 34 bp to 3303 bp; peptide length: 1090 Category: similarity to unknown protein Classification: no clue

1 MGIQNLQRSE SSKMDKYETE ESSVAGLSSP ELKVRPAGAS SIWYTEGEKQ 51 LTKSLKGKNE ESNKSKVKVT KLMKTMKSEN TKKLIKQNSK DSVVLVGYKC JOB LKSTASNDLI KCFEGNPSHS QKEGLDPTIC GYNFDPKTYM RQTSQKEASC 35 151 LPTNTERTEQ KSPDIENVQP DQFDPLNSGN LNLCANLSIS GKLDISQDDS 201 EITQMEHNLA SRRSSDDCHD HQTTPSLGVR TIEIKPSNKD PFSGENITVK 251 LGPWTELRQE EILVDNLLPN FESLESNGKS KSIEITFEKE ALQEAKCLSI 301 GESLTKLRSN LPAPSTKEYH VVVSGDTIKL PDISATGRASH NZRJSSSESE 351 PSSFATHPNT DLVFETVQGQ GPCNSERLFP QLLMKPDYNV KFSLGNHCTE 40 407 ZIZYIZEIGZ ZTIZINZTЬZ DDETZЬDENZ KKZANЬECHT NDZKIATUT 451 TTDLPKCDDT KKSSITLQQQ SVVFSGNLDN ETVAIHSLNS SIKDPLQFVF 501 SDEETSSDVK SSCSSKPNLD TMCKGFQSPD KSNNSTGTAI TLNSKLICLG 551 TPCVISGSIS SNTDVSEDRT MKKNSDVLNL TQMYSEIPTV ESETHLGTSD PESASIDIAK GELAENALES GSSIDISDIC AASASUMTSE GKEIZEKEIS 651 NLQQEQDKED EEEEQDQQMV QNGYYEETDY SALDGTINAH YTSRDELMEE 701 RLTKSEKINS DYLRDGINMP TVCTSGCLSF PSAPRESPCN VKYSSKSKFD 751 AITKQPSSTS YNFTSSISWY ESSPKPQIQA FLQAKEELKL LKLPGFMYSE BOD VPLLASSVPY FSVEEEGGSE DGVHLIVCVH GLDGNSADLR LVKTYIELGL 851 PGGRIDFLMS ERNQNDTFAD FDSMTDRLLD EIIQYIQIYS LTVSKISFIG 50 POI HSLGNLIIRS VLTRPRFKYY LNKLHTFLSL SGPHLGTLYN SSALVNTGLW 951 FMQKWKKSGS LLQLTCRDHS DPRQTFLYKL SNKAGLHYFK NVVLVGSLQD 1001 RYVPYHSARI EMCKTALKDK QSGQIYSEMI HNLLRPVLQS KDCNLVRYNV

1051 INALPHTADS LIGRAAHIAV LDSEIFLEKF FLVAALKYFQ

55

BLASTP hits

No BLASTP hits available

Alert BLASTP hits for DKFZphtes3_26g3, frame 1

5 No Alert BLASTP hits found

EKWI

SEQ

SEG

Pedant information for DKFZphtes3_2bg3, frame 1

10 Report for DKFZphtes3_26g3.1

LOW_COMPLEXITY

1701 **ELENGTHI** 155542.55. EMMI 15 [pI] 5-12 EHOMOLE TREMBL:CEAF2196_1 gene: "CO9D4.4"; Caenorhabditis elegans cosmid CO9D4. 2e-38 EFUNCATI 99 unclassified proteins ES. cerevisiae, Y0RD59cJ 2e-06 BF00750B 20 EBF0CK21 EKWI Alpha_Beta

6.72 %

DSVTEDLDAPWMGIQNLQRSESSKMDKYETEESSVAGLSSPELKVRPAGASSIWYTEGEK . 25 SEQ SEG PRD ccccccccceeeechhhhhhhhhhhccccccccccccceeeeccccch **GLTKZLKGKNEESNKZKVKVTKLMKTMKSENTKKLIKQNZKDSVVLVGYKCLKZTAZNDL** SEQ 30 SEG -----xxxxxxxxxxxxxxxxxx PRD IKCFEGNPSHSQKEGLDPTICGYNFDPKTYMRQTSQKEASCLPTNTERTEQKSPDIENVQ SEQ SEG 35 SEQ : PDQFDPLNSGNLNLCANLSISGKLDISQDDSEITQMEHNLASRRSSDDCHDHQTTPSLGV SEG 40 SEQ : RTIEIKPSNKDPFSGENITVKLGPWTELRQEEILVDNLLPNFESLESNGKSKSIEITFEK SEG PRD EALQEAKCLSIGESLTKLRSNLPAPSTKEYHVVVSGDTIKLPDISATYASSRFSDSGVES 45 SEQ SEG PRD EPSSFATHPNTDLVFETVQGQGPCNSERLFPQLLMKPDYNVKFSLGNHCTESTSAISEIQ ZEQ 50 SEG PRD ZSLTSINSLPSDDELSPDENSKKSVVPECHLNDSKTVLNLGTTDLPKCDDTKKSSITL@@ SEQ SEG 55 PRD

@ZVVFZGNLDNETVAIHZLNZSIKDPL@FVFZDEETZZDVKZZCZZKPNLDTMCKGF@ZP

_____xxxxxxxxxxxxx

	PRD	eeeeeeccccceeeeeeccccceeeeeccccceeeecccc
5	SEQ SEG PRD	DKZNNSTGTAITLNSKLICLGTPCVISGSISSNTDVSEDRTMKKNSDVLNLT@MYSEIPT
10	SEQ SEG PRD	CCCCCCCCCCCCCCGGGGGGGGGGGGGGGGGGGGGGGG
10	SEQ SEG PRD	SNLQQEQDKEDEEEEQDQQMVQNGYYEETDYSALDGTINAHYTSRDELMEERLTKSEKINxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
15	SEQ SEG PRD	ZDYLRDGINMPTVCTSGCLSFPSAPRESPCNVKYSSKSKFDAITKQPSSTSYNFTSSISW
20	SEQ SEG PRD	CCCCCCCCHUHHHHHHHHHHHCCCCCeeeeeeeeeeeccccceeeeeee
25	SEQ SEG PRD	HGLDGNSADLRLVKTYIELGLPGGRIDFLMSERNANDTFADFDSMTDRLLDEIIAYIAIY ecccccchhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
30	SEQ SEG PRD	SLTVSKISFIGHSLGNLIIRSVLTRPRFKYYLNKLHTFLSLSGPHLGTLYNSSALVNTGL hcccccccccccceeeeecccccccccccccccccccc
30	SEQ SEG PRD	WFMQKWKKSGSLLQLTCRDHSDPRQTFLYKLSNKAGLHYFKNVVLVGSLQDRYVPYHSAR
35	SEQ- SEG PRD	IEMCKTALKDKQSGQIYSEMIHNLLRPVLQSKDCNLVRYNVINALPNTADSLIGRAAHIA
40	SEQ SEG PRD	VLDSEIFLEKFFLVAALKYFQ hhhhhhhhhhhhhhhhccc
	(No	Prosite data available for DKFZphtes3_2bg3.1)
45	(No	Pfam data available for DKFZphtes3_26g3.1)

DKFZphtes3_29f24

5 group: signal transduction

DKFZphtes3_29f24 encodes a novel 526 amino acid protein with similarity to murine netla.

10 The closely related mNETL activates signalling pathways in addition to those directly controlled by activated RhoA. The novel protein is expressed ubiquitously.

The new protein can find application in modulation/blocking signalling pathways.

similarity to netla (Mus musculus)

20 perhaps complete cds.

Sequenced by BMFZ

Locus: /map="72.40 cR from top of Chr3 linkage group"

25 Insert length: 3559 bp

Poly A stretch at pos. 3534, polyadenylation signal at pos. 3513

1 CGCCGCCGC CGGCATCGTG GAGCTGGGGC CCCCTTTTGC CTGGGAGTTT 30 51 TGTAGTCGCC TAGGGTCAGC GGTGACATCC CAAAGGGCAG GCCCGGCAGC 151 ACTGCAGCCT GGAGCTACCC CCGGCCAGCG GTCCGGCCAA GGACGCTGAG 35 451 TTACATCCAA GGAAATCAAA CGTCAGGAGG CGATCTTTGA GCTTTCCCAA 503 GGAGAAGAAG ACTTGATAGA AGACTTGAAA TTAGCAAAAAA AGGCCTATCA
553 TGACCCCATG CTGAAACTCT CCATAATGAC AGAACAAGAG TTGAATCAAA
603 TTTTTTGGAAC ACTGGACTCT CTAATTCCTC TACATGAAGA GCTCCTTAGT
653 CAGCTTCGAG ATGTTAGGAA GCCTGATGGC TCGACTGAAC ATGTTGGTCC
703 CATCCTCGTG GGCTGGCTCC CTTGCCTCAG CTCCTATGAT AGCTACTGCA
751 GCAATCAAGT AGCCGCCAAA GCTCTTTA CAATCCCCCAT 40 45 BDL CGAGTCCAGG ATTTCCTACA GCGATGTTTA GAATCCCCCT TTAGCCGCAA 851 ACTAGATETE TEGAATTTEE TEGATATTEE AAGAAGEEGE CTEGTAAAAT POL ACCCTCTGCT TCTCCGAGAA ATCTTGAGGC ACACACCAAA TGATAATCCA 951 GATCAGCAGC ACTTGGAAGA AGCTATAAAT ATCATTCAGG GAATTGTGGC DOD AGAAATCAAC ACCAAGACTG GTGAATCTGA ATGCCGCTAT TATAAAGAGC 50 1051 GGCTTCTTTA CTTGGAAGAA GGCCAGAAAG ACTCCCTGAT CGACAGCTCT ኔኒበኔ CGAGTCTTGT GTTGTCATGG TGAACTGAAG AACAATCGGG GCGTGAAACT 1151 GCATGTTTTC CTGTTCCAAG AAGTGCTTGT GATCACTCGA GCCGTCACCC 1201 ACAATGAGCA GCTTTGCTAC CAGCTGTACC GTCAGCCAAT CCCCGTGAAA 1251 GACCTCCTGC TGGAAGACCT CCAGGATGGA GAAGTGAGGC TGGGTGGCTC 55 LBOL CCTGCGAGGG GCATTCAGCA ACAATGAGAG AATTAAAAAC TTCTTCAGAG 1351 TCAGTTTCAA AAATGGATCC CAAAGTCAGA CCCACTCGCT ACAAGCCAAT 1401 GACACTTTCA ACAAACAGCA GTGGCTTAAC TGTATTCGTC AAGCCAAAGA

	1451	AACAGTTTTG	TGTGCTGCCG	GGCAAGCTGG	GGTGCTTGAC	TCCGAGGGAT
	1501	CGTTCCTAAA	TCCCACCACC	GGGAGCAGAG	AGCTACAGGG	AGAAACAAAA
	1551	CTTGAGCAGA	TGGACCAATC	GGACAGTGAG	TCAGACTGTA	GTATGGACAC
	1601	GAGTGAGGTC	AGCCTCGACT	GTGAGCGCAT	GGAACAGACA	GACTCTTCCT
5	1651	GTGGAAACAG	CAGGCACGGT	GAAAGTAACG	TCTGACAGAA	GCATGTGCAC
_	1701	TTCGGGAAGC	AGGCCTGCAT	CTTACCTGTA	CAGTATTTGC	ATTCCACAGA
	1751	TGGAACGGTT	TGGAGAAGCA	CTTTTTCATA	CTTTTGTGAA	AGTATACATG
	1801	TTGGCCCAGT	CTCTCGTATC	TGTACCTTTG	TCCCTAGTAC	TGTAACTGCC
	1851	AATCTGTCTG	TGTAAGCTGG	AATCTGTGGC	AACTATTACC	CTGTGTTGTA
10	1901	TTTCCCAAGT	GTCTGGATGG	ATGGAGAGGT	ACTCAAACAA	GTTACTTTCA
10	1951	GTTGTCCTGC	TGGATTTTAA	AAAATAGAA	AAAGAATCTC	AAAACTACTG
	5007	TTTTACATAG	ATTGTTTGAA	GAGTCCTTCC	TCTTGTGCTT	CTGTACCACT
	2051	TTCCCAGCTC	TTAGATGTGG	TAGCTAAAGG	CACGGAATTT	AGACGGCCTT
	5707	GTAAATAGGG	CATGAGGAAC	TCATCTGTGT	ATTGGGATGG	TATTAGAGAG
15	2151	AGAATCAGGA	AAGACCAACT	CATGAAGTGA	ACTTGGTTTG	ATCTTACTCA
15	5507	ACTAGAAAGC	TTGAAAACAT	CCCTGGGGAT	TCTGAAGGCT	TAATTTTGCA
	2251	AAGGAGGATG	CATTGTCTGA	ACTTTGCAAC	TTCATCCAGT	GCAAGTTTGA
	5307	TGCAAGAATG	TATTAGGACA	TAAAATAGAG	GCTGACCTTA	AAAGGCCAG
	2351	GACAGAAGCG	GCTGCCAGCT	CTGAATCTTT	AACTGAAATG	CACATGGCAC
20	2407	CAGGAGGTGT	CTCTCATAGT	TGGTTGCTAG	CCTAAAACAT	CAGAATAGAA
20	2451	CCCAAAGGGC	TTAGGAAGGC	CTGCCAGGAT	AACAAGAAGG	CCCTGTATTC
	2501	ATTGTGTTTC	ATCTGCCTAG	GCCTACTCAT	TATTTTAGAG	AATGAATGAA
	2551	GCAACAAGGA	AGAGAGACCA	TGACTCTATC	GATGACACTG	TTTATAGAAA
	5207	CACAGGAGAG	GAAGAATTTG	GAATGAAAAG	CACTTCGTCA	GAACCTTCTG
25	2651	TGGGAGCCAT	TGAGAGAAAA	GCATGGTCCA	GTGCCTTCTG	AGAAAGGCCA
LJ	2701	GAGCTTTGGG	CTTTCCTGCT	CTGCTTTTGG	GTCGTCAATT	TGCCATCTCT
	2751	GGTTCTGTGC	TATAATCAGA	ATTGTAATTA	TGTTCTCCAG	AGGCCAATTT
	5907	CATTAACTCT	GATTAATTAG	AATCAGCTAG	CCAGATTAGT	AACCTCTTTG
	2851	TCCAGCCTTG	ATTTACAGTG	CAGGGTAAAG	TGCAGACCTT	AAAAACAGCT
30	5907	AAGTACCTAG	AAGAGCTCCC	TGCAAGTGTA	AATATTAAGG	ATGACCTGTG
20	2951	CAAAATTATA	CCCACACCAG	CACTAGTGGT	AATTATTCTA	AATTATTGCC
	3001	AAAAAGTTTT	TTTTAATCTG	TCTTTCAAGT	TTACAGAAAA	GAAAGCAGTA
	3051	AATGCATTGA	TGTCATTTTA	TTATGTACAT	ATATCATGTG	CATTCAAGCT
	. 3101	GTGTGACAAG	ATATATCAAT	ATAAAAACAA	GGTATATACT	TTATTATTTT
- 35		-TTGAAAACAA		ATCAATTTTA	CCCTGTAAAA	.CATATTTCT.G.
•••		TATTTATAGG	TCTTAAACAT	GATGAATTTT	TTCTATTACA	AGTTTATTTA
	3251	AAACTGCTTT	CTCAAGTCGT	TATTGATACA	GCAAGTGAAC	CTGCTGCAGA
	3301	CAGAAGCAGA	GGAAAGCCAA	GAACAGCCTT	TATTGGTGAA	GAAAAGAATG
	3351	AATGATTCTT	TGTAGGCGCC	ATCAGCCACT	TTTAGAAGCC	ATCAGCCAGT
40	3401	GTGTTGGGAA	AAGAGGTTTG	TCAAGTGTTG	GCCTATGGGA	AGGTGGTCAA
	3451	TGAATGTTTT	GATGAAATGA	ATGTTTTTGT	ATAATGGCCT	TAAACTTTTC
	3501	TGGAAGTATT	TCAAATAAAT	TACATTATTA	AGTCAAAAA	AAAAAAAAA
	3551	AAAAAAAA				

45

BLAST Results

No BLAST result

50

Medline entries

55 98336196:
Alberts AS: Treisman R: Activation of RhoA and SAPK/JNK signalling pathways by the

WO 01/98454



RhoA-specific exchange factor mNETL. EMBO J 1998 Jul 15:17(14):4075-85

5

Peptide information for frame 3

10 ORF from 105 bp to 1682 bp; peptide length: 526 Category: strong similarity to known protein Classification: Cell signaling/communication

1 MVAKDYPFYL TVKRANCSLE LPPASGPAKD AEEPSNKRVK PLSRVTSLAN 51 LIPPVKATPL KRFSQTLQRS ISFRSESRPD ILAPRPWSRN AAPSSTKRRD 15 JOJ SKLWSETFDV CVNQMLTSKE IKRQEAIFEL SQGEEDLIED LKLAKKAYHD 151 PMLKLSIMTE QELNQIFGTL DSLIPLHEEL LSQLRDVRKP DGSTEHVGPI 507 FARMEDCEZZ ADZACZNGAV VKVFFDHKKG DHKAGDŁFGK CFEZЬŁZKF 251 DLWNFLDIPR SRLVKYPLLL REILRHTPND NPDQQHLEEA INIIQGIVAE .20 HANDANNAH BOHDANASZ DISTANTIN EEGQKDSLID SSRVLCCHGE LKNNAGYKLH 351 VFLFQEVLVI TRAVTHNEQL CYQLYRQPIP VKDLLLEDLQ DGEVRLGGSL 401 RGAFSNNERI KNFFRVSFKN GSQSQTHSLQ ANDTFNKQQW LNCIRQAKET 451 VLCAAGQAGV LDSEGSFLNP TTGSRELQGE TKLEQMDQSD SESDCSMDTS 501 EVSLDCERME QTDSSCGNSR HGESNV

25

BLASTP hits

. 30 No BLASTP hits available

Alert BLASTP hits for DKFZphtes3_29f24, frame 3

No Alert BLASTP hits found

35

Pedant information for DKFZphtes3_29f24, frame 3

Report for DKFZphtes3_29f24.3

40

50

560 **ELENGTHI** EMMI P3505·82 6.04 [[q]

TREMBL:AF094520_1 gene: "Netl"; product: "NETL 45 [HOMOL] homolog"; Mus musculus NETL homolog (Netl) mRNA, complete cds. le-195 **EFUNCATE** 09.01 biogenesis of cell wall ES. cerevisiae,

YLR371w1 3e-16

EFUNCATE 03.07 pheromone response, mating-type determination, **EFUNCATE 10.02.09** regulation of g-protein activity

cerevisiae, YLR371w1 Je-16 EFUNCATI 09.04 biogenesis of cytoskeleton ES. cerevisiae.

55 YLR371w3 3e-16

EFUNCATE 03.04 budding, cell polarity and filament formation ES. cerevisiae, YLR371wl 3e-16

WO 01/98454 PCT/IB01/02050 EFUNCATD 01.05.04 regulation of carbohydrate utilization cerevisiae, YLR371w1 3e-16 EFUNCATD 30.03 organization of cytoplasm IS. cerevisiae. YALD41w1 3e-11 EFUNCATI 03.22 cell cycle control and mitosis 5 IIS. cerevisiae, YAL041w1 3e-11 EFUNCATD 10.05.09 regulation of g-protein activity EZcerevisiae, YALO41wl 3e-11 EBF0CKZ] PR00510E 10 **EBFOCKZ** PRODDUJE EBF0CK21 BL00741B EPIRKW1 breakpoint cluster region le-Db **EPIRKUD** transmembrane protein 5e-13 **EPIRKU**I brain 3e-Ob 15 **EPIRKU** signal transduction 5e-13 **EPIRKU**J alternative splicing le-Ob **ESUPFAMD** CDC24 homology 9e-15 ESUPFAMI SH2 homology le-ll **ISUPFAMD** CDC25-type quanine nucleotide exchange activator 20 homology 2e-08 ESUPFAMI dbl transforming protein 9e-08 protein kinase C zinc-binding repeat homology le-ll [SUPFAM] ESUPFAMI SH3 homology le-ll **ESUPFAMD** bcr protein le-Ob 25 **ESUPFAM3** pleckstrin repeat homology 2e-11 **CSUPFAMJ** vav transforming protein le-ll EKWI All_Alpha 30 SEQ PPPGIVELGPPFAWEFCSRLGSAVTSQRAGPAAMVAKDYPFYLTVKRANCSLELPPASG PRD SEQ PAKDAEEPSNKRVKPLSRVTSLANLIPPVKATPLKRFSQTLQRSISFRSESRPDTI APRP PRD 35 **USRNAAPSSTKRRDSKLWSETFDVCVNQMLTSKEIKRQEAIFELSQGEEDLIEDLKLAKK** SEQ PRD SEQ AYHDPMLKLSIMTEQELNQIFGTLDSLIPLHEELLSQLRDVRKPDGSTEHVGPILVGWLP 40 PRD SEQ CLSSYDSYCSNQVAAKALLDHKKQDHRVQDFLQRCLESPFSRKLDLWNFLDIPRSRLVKY PRD 45 SEQ PLLLREILRHTPNDNPDQQHLEEAINIIQGIVAEINTKTGESECRYYKERLLYLEFGQKD PRD SLIDSSRVLCCHGELKNNRGVKLHVFLFQEVLVITRAVTHNEQLCYQLYRQPIPVKDIII SEQ PRD 50 EDLQDGEVRLGGSLRGAFSNNERIKNFFRVSFKNGSQSQTHSLQANDTFNKQQWLNCIRQ SEQ PRD AKETVL CAAGQAGVLDSEGSFLNPTTGSRELQGETKLEQMDQSDSESDCSMDTSEVSLDC SEQ PRD 55

-366-

SEQ

PRD

ERMEQTDSSCGNSRHGESNV

cccccccccccccc

(No Prosite data available for DKFZphtes3_29f24.3)

5 (No Pfam data available for DKFZphtes3_29f24.3)

DKFZphtes3_30pb

5 group: testis derived

DKFZphtes3_30pb encodes a novel 4bl amino acid protein without similarity to known proteins.

10 No informative BLAST results; No predictive prosite, pfam or SCOP motife.

The new protein can find application in studying the expression profile of testis-specific genes.

15

similarity to C.elegans F41H10-4

perhaps complete cds.

20

Sequenced by LMU

Locus: unknown

25 Insert length: 1944 bp

Poly A stretch at pos. 1911, no polyadenylation signal found

1 GGAACAGACC ACTGGGCTGG CAGCTGAGTT GCAGCAGCAG CAGGCTGAGT 30 51 ACGAGGACCT TATGGGACAG AAAGATGACC TCAACTCCCA GCTCCAGGAG LOL TCATTACGGG CCAATAGTCG ACTGCTGGAA CAACTTCAAG AAATAGGGCA 151 GGAGAAGGAG CAGTTGACCC AGGAATTACA GGAGGCTCGG AAGAGTGCGG 201 AGAAGCGGAA GGCCATGCTG GATGAGCTAG CAATGGAAAC GCTGCAAGAG 251 AAGTCCCAGC ACAAGGAAGA GCTGGGAGCA GTTCGTCTAC GGCATGAGAA 3D1 GGAGGTGCTG GGGGTGCGTG CCCGCTATGA GCGTGAGCTC CGAGAGCTGC 35 351 ATGAAGACAA GAAGCGTCAG GAGGAGGAC TCCGTGGCA GATCCGGGAG 401 GAGAAGGCCC GGACACGGGA GCTGGAGACT CTCCAGCAGA CAGTGGAAGA 451 ACTTCAAGCT CAGGTACATT CCATGGATGG AGCCAAGGGC TGGTTTGAAC 5D1 GGCGCTTGAA GGAAGCCGAG GAATCCCTGC AGCAGCAGCA GCAGGAACAA 551 GAGGAAGCCC TCAAGCAGTG TCGGGAGCAG CACGCTGCCG AGCTGAAGGG 40 LOD CAAGGAGGAG GAGCTACAGG ATGTACGGGA TCAGCTCGAG CAGGCCCAGG **L51 AGGAGCGGGA CTGCCACCTG AAGACCATTA GCAGCCTGAA GCAGGAGGTG** 7D1 AAGGACACAG TGGATGGGCA GAGGATCCTG GAGAAGAAGG GCAGTGCTGC 751 GCTCAAGGAC CTCAAGCGGC AGCTGCATTT GGAGCGGAAA CGGGCAGATA 45 BD1 AGCTGCAGGA GCGACTGCAG GACATCCTCA CTAACAGCAA GAGCCGCTCA 851 GGCCTTGAGG AGCTGGTTCT CTCAGAGATG AACTCACCAA GCCGGACCCA 9D3 GACAGGGGAC AGCAGTAGCA TCTCCTCCTT CAGCTACCGG GAGATCTTGC 951 GGGAAAAGGA GAGCTCGGCT GTTCCAGCCA GGTCCTTATC CAGCAGCCCT LODL CAAGCCCAGC CCCCTCGGCC AGCAGAGCTG TCAGATGAGG AAGTGGCTGA 1051 GCTCTTTCAG CGGCTGGCAG AGACACAGCA GGAGAAATGG ATGCTGGAGG 50 . LIDL AGAAGGTGAA GCACCTGGAA GTGAGCAGTG CTTCCATGGC AGAGGACCTC 1151 TGCCGGAAGA GCGCCATCAT TGAGACCTAC GTCATGGACA GCCGGATCGA 1201 TGTGTCTGTG GCAGCAGGCC ACACAGACCG CAGCGGGCTG GGCAGCGTCC 1251 TGAGAGACCT AGTGAAGCCA GGCGACGAGA ACCTTCGGGA GATGAACAAG ACOTTOADOT ATAADAADOA OTOBADDA ADATOATAA ADACOTOAA LOEL 55 1351 CAAGGATATG GAAGTTCTGT CCCAGGAAAT TGTGCGGCTC AGCAAGGAGT
1401 GCGTGGGGCC TCCTGACCCA GACCTAGAGC CAGGAGAAAC CAGCTAAAGA
1451 CCTGCAGGCT GCACCCACCT CCTCCCCTTC CTACCCCCTA GGATGCTATT

WO 01/98454 PCT/IB01/02050 1501 CCCTTGGGCT GTGGTGGAAA AATGAGGGCT GGAGCCAAAA TCAAATAGCT 10 **BLAST Results** 15 No BLAST result Medline entries ______ 20 No Medline entry 25 Peptide information for frame 2 ORF from 62 bp to 1444 bp; peptide length: 461 Category: similarity to unknown protein 30 Classification: no clue 1 MGQKDDLNSQ LQESLRANSR LLEQLQEIGQ EKEQLTQELQ EARKSAEKRK 51 AMLDELAMET LØEKSØHKEE LGAVRLRHEK EVLGVRARYE RELRELHEDK 101 KRQEEELRGQ IREEKARTRE LETLQQTVEE LQAQVHSMDG AKGWFERRLK 151 EAEESLAARA REGEEALKAC REGHAAELKG KEEELADVRD QLEGAGEERD 35 201 CHLKTISSLK QEVKDTVDGQ RILEKKGSAA LKDLKRQLHL ERKRADKLQE 251 RLQDILTNSK SRSGLEELVL SEMSPSRTQ TGDSSSISSF SYREILREKE 301 SSAVPARSLS SSPQAQPRP AELSDEVAE LFQRLAETQQ EKWMLEEKVK JURIVZDIDZ ROTHDAAVZV DIRZDMVYTE IIAZNROLDE AMZAZZVEJL 40. 4D1 VKPGDENLRE MNKKLQNMLE EQLTKNMHLH KDMEVLSQEI VRLSKECVGP 451 PDPDLEPGET S 45 BLASTP hits No BLASTP hits available Alert BLASTP hits for DKFZphtes3_30pb, frame 2 50 No Alert BLASTP hits found Pedant information for DKFZphtes3_30pb, frame 2 55

Report for DKFZphtes3_30p6.2

WO 01/98454 PCT/IB01/02050

ELENGTHD 481

	ELENGTHD 481	
	EMW] 55398·10	
	[pI] 5.0?	
5	<pre>LHOMOLD TREMBL:CEF41H10_4 gene: "F41H10.4"; Caenorhabditi: elegans cosmid F41H10. 2e-12</pre>	S
_	EFUNCATE 30.03 organization of cytoplasm ES. cerevisiae	
	YDL058w3 5e-04	
	EFUNCATD O8.07 vesicular transport (golgi network, etc.) ES.	
	cerevisiae YDLO58w3 5e-04	
10	EBLOCKSD BLOLLOOD NNMT/PNMT/TEMT family of methyltransferases	
	proteins	
	ĽΚຟƊ All_Alpha	
	EKWD LOW_COMPLEXITY 19.13 %	
	EKWD COILED_COIL 40-96 %	
15	•	
	SEQ EQTTGLAAELQQQQAEYEDLMGQKDDLNSQLQESLRANSRLLEQLQEIGQEKEQLTQEL	^
	SEQ EQTTGLAAELQQQQAEYEDLMGQKDDLNSQLQESLRANSRLLEQLQEIGQEKEQLTQEL SEG×××××××××××××××××××××××××××××××××	
	PRD ccchhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh	
20	COILS	1 8
		c
		_
	SEQ EARKSAEKRKAMLDELAMETLQEKSQHKEELGAVRLRHEKEVLGVRARYERELRELHEDI	K
	SEG x	•
25	PRD հորհիրիկիրիկիրիկիրիկիրիկիրիկիրիկիրիկիրիկիր	h
	COILS	
	ccccccc	•
	SEQ KRQEEELRGQIREEKARTRELETLQQTVEELQAQVHSMDGAKGWFERRLKEAEESLQQQ	^
30	ZEGXXXXXXXXXXXXXXXXXXXXXXIIIQAKGWFEKKKEAEEZERKK	
50	PRD hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh	
	COILS	•
		C
	•	
35	SEQ QEQEEALKQCREQHAAELKGKEEELQDVRDQLEQAQEERDCHLKTISSLKQEVKDTVDG	Q
	SEG xxxxxxxx	-
	PRD հիրանի հերանական անագրագրերը հետում է հերանական հետում է հետում է հետում է հետում է հետում է հետում է հետու	C
	CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC	
40		•
70	SEQ RILEKKGSAALKDLKRQLHLERKRADKLQERLQDILTNSKSRSGLEELVLSENNSPSRT	^
	SEG	×
	PRD cccccchhhhhhhhhhhhhhhhhhhhhhhhhhccccchhhh	_
	COILS	_
45		
	SEQ TGDSSSISSFSYREILREKESSAVPARSLSSSPQAQPPRPAELSDEEVAELFQRLAETQ	
	SEG ···××××××××······××××××××××××××××××××	
~ 0	PRD cccccchhhhhhhhhhhhccccccccccccccchhhhhh	h
50	COILZ	
	•••••••••••••••••••••••••••••••••••••••	•
	CEN EVUMI EEVUVUI EUGGAGMAERI ZEVGATTETVUMRGETRUGULAZUTREGG GOV GOV	
	SEQ EKWMLEEKVKHLEVSSASMAEDLCRKSAIIETYVMDSRIDVSVAAGHTDRSGLGSVLRD SEG	-
55	PRD hhhhhhhhhhhhhhhhhhhhhhhhhhccccchhhhhhhh	•
	Colls	-

	W	O 01/98454	PCT/IB01/02050		
	SEQ	VKPGDENLREMNKKLØNMLEEQLTKNMHLHKDMEVLSØEIVRLSK	ECVGPPDPDLEPGET		
	SEG				
	PRD	cccchhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh	ccccccccccc		
	COIL	2			
5		cccccccccccccccccccccccccccccccccc			
	SEQ	2			
	SEG	•			
	PRD	с			
10	COIL	z			
	(No	Prosite data available for DKFZphtes3_30p6.2)			
15	(No	Pfam data available for DKFZphtes3 30pb.2)			

DKFZphtes3_3lalO

5 group: nucleic acid management

DKFZphtes3_3lalO encodes a novel 542 amino acid protein with similarity to histone Hl of Drosophila hydei.

10 Histone HL variants are known to act as specific regulators of genes via the differential condensation of DNA.

The new protein can find application in modulating/blocking the transcriptional activity and in expression profiling.

15

weak similarity to Drosophila histone Hl

perhaps complete cds.

20

Sequenced by LMU

Locus: /map="13"

25 Insert length: 2887 bp
Poly A stretch at pos. 2855, polyadenylation signal at pos. 2839

B AGATGATCCC CAAAGTCAAC ATATGACATT AAGCCAGGCA TTTCACCTTA 51 AAAACAATAG TAAAAAGAAA CAAATGACTA CAGAAAAACA AAAGCAAGAT 30 DDL GCTAACATGC CCAAGAAACC TGTGCTTGGA TCTTATCGTG GCCAGATTGT 151 TCAGTCTAAG ATTAATTCAT TTAGAAAACC TCTACAAGTC AAAGATGAGA 201 GTTCTGCAGC AACAAAGAAA CTTTCAGCCA CTATACCTAA AGCCACAAAA 251 CCTCAGCCTG TAAACACCAG CAGTGTAACA GTGAAAAGTA ATAGATCCTC 351 AACTTGTGCG ACCTCCTATT AGAAGTCATC ACAGTAATAC CCGGGACACT 401 GTGAAACAAG GCATCAGTAG AACCTCTGCC AATGTTACAA TCCGGAAAGG 451 GCCTCATGAA AAAGAACTAT TACAATCAAA AACAGCTTTA TCTAGTGTCA 501 AAACCAGTTC TTCTCAAGGT ATAATAAGAA ATAAGACTCT ATCAAGATCC 551 ATAGCATCTG AAGTTGTAGC CAGGCCTGCT TCATTGTCTA ATGATAAACT 40 LOL GATGGAAAAG TCAGAGCCCG TTGACCAGCG AAGACATACT GCAGGAAAAG LSI CAATTGTTGA TAGTAGATCA GCTCAGCCCA AAGAAACCTC GGAAGAGAGA 701 AAAGCTCGTC TGAGTGAGTG GAAAGCTGGC AAAGGAAGAG TGCTAAAAAG 751 GCCCCCTAAT TCAGTAGTTA CTCAGCATGA GCCTGCAGGA CAAAATGAAA BOL AACTAGTTGG GTCTTTTTGG ACTACCATGG CAGAAGAAGA TGAACAAAGA 45 B51 TTATTTACTG AAAAAGTAAA CAACACATTT TCTGAATGCC TGAACTTGAT POL TAATGAGGGA TGTCCAAAAG AAGATATACT GGTCACACTG AATGACCTGA 951 TTAAAAATAT TCCAGATGCC AAAAAGCTTG TTAAGTATTG GATATGTCTT TTAAAAATAT TCCAGATGCC AAAAAGCTTG TTAAGTATTG GATATGTCTT

LOOL GCACTTATTG AACCAATCAC AAGTCCTATT GAAAATATTA TTGCAATCTA

LOSL TGAGAAAGCC ATTCTGGCAG GGGCTCAGCC TATTGAAGAG ATGCGACACA

LLOL CGATTGTAGA TATTCTAACA ATGAAGAGTC AAGAAAAAGC TAATTTAGGA

LSL GAAAATATGG AGAAGTCTTG TGCAAGCAAG GAAGAAGAGTCA AAGAAGTCAG

LSOL TATTGAAGAT ACAGGTGTTG ATGTAGATCC AGAAAAAACTG GAAATGGAGA

LSSL GTAAACTTCA TAGAAATTTG CTATTTCAAG ATTGTGAAAA AGAGCAAGAC

LSOL AACAAAACAA AAGATCCAAC CCATGATGTT AAAACCCCCA ATACAGAAAC

LSSL GAGGACAAGT TGCTTAATTA AATATAATGT GTCTACTACG CCATACTTGC

LHOL AAAGTGTGAA AAAAAAGGTG CAGTTTGATG GAACAAATTC CGCATTTAAA

LHSL GAGCTGAAGT TTTTAACACC AGTGAGACGT TCTCGACGTC TTCAAGAGAA 50 55

WO 01/98454 PCT/IB01/02050 WO 01/98454

DESID AACTTCTAAA TTGCCAGATA TGTTAAAAGA TCATTATCCT TGTGTGTCTT

DESID CATTGGAACA GCTAACGGAG TTGGGAAGAG AAACTGATGC TTTTGTATGC

LLDD CGCCCTAATG CAGCACTGTG CCGGGTGTAC TATGAGGCTG ATACAACATA

LLSD AGAGAAATAA AGCTCTGTTA GGGAATGGGG TTTTTATTAT TTGTGGGGTG

DTD TTTTGTTTTG AGTAGCTTTA TATTGCTCTT AGGTCTGGAG TTGGCCATGT

DTD TTTTGTTTTG AGTAGCTTTA TATTGCTCTT AGGTCTGGAG TTGGCCATGT

DTD TGTTATGGCA AGAGTTGTCC TCTACATTGG AAAGCTAATC CTACCTTGTC

DBD TGTTATGGCA AGAGTTGTCC TCTACATTGG AAAGCTAATC CTACCTTGTC

DBD AGTTTACTAT GTTCCTTGAA TATAAACAGG TTATAATACT ACCCTGTTCA

DTD AGTTTACTAAA TATAAGTACA GTAATGATGC ATAATTAGAA AATGAGGTAT

DDD TCTAGGTAAA ATGTATGTTT GCCTTGACAT GTTTTTAAAA GTTATGATGT

DDD AGATTAGTCA AAAATTCTAT AGAATGACTC ACTTCGAATA CTAAGACACA

DDD AGATTAGTCA AAAATTCTAT AGAATGACTC ACTTCGAATA CTAAGACACA

DDD CTCTTGCGTC CCTTGGACTG CCTGTTGATT GATGTAACCT

DCTCTTGCGTC CCTTGGACTG CCTGTTGATT GATGGAAAGT GTCTGCACTG

DCTCTTGCGTC CCTTGGACTG CCTGTTTGATT GATGGAAAGT GTCTGCACTG

DCTCTTGCACCA AGAAGGTTTA CTTAAATTAAA TACCGCATTT CTAAGAGAAAG 5 10 15 23D1 GGGGTCACCA AGAAGGTTTA CTTAATTAAA TACCGCATTT CTAAGAGAAG 2351 ATACTTTGT TAAGAAAGA TGCCACATTT AGTGGTTTAA CTTTTGTAAC 2451 TCACTGAT AGTTTTTAAG CAATTAGAT GGAGTTAGG AAAGAACATA 2451 TCACACACA CAAATGTCAT TCTAGTTAAG ATAGCATTC TAAGATAACT 20 2501 GATACTGAA CAAATGTCAT TCTAGTTTAG ATAGCATTTC TAAGATAACT
2501 GATACTAATA CTTGTTTTCT TCCCTATAAC ATAAAAAACT TCACTGTTAA
2551 GTCATGTCCC TTGAAACATG ATAGTTACAT ACACAGTTTT CTCTCCACAC
2601 ATAAATAACA CCACTAAAGT TGTTTTGTAA GGTTCCAAAC TAATATGGCA
2651 TATATCAACT CTACAGTTTC AAATAAATGA CTTTTTAATT GTAAAAGATT
2701 AGTTGAAAAA CTGTATGAAT GTGAAGATCA CATGCTTAGT CATTTTTATG
2751 TTCATTCCAC TTTGTATATC TTTTCTATTT ATTGACTTCT CATGTTCTAG
2801 AGAGTAGGAC TTTTATTCCG TGTACCTGAT ATATATACAA TTAAAATATC 25 2851 TGTATAATTA AAAAAAAAA AAAAAAAAA AAAAAAG 30 **BLAST Results** No BLAST result 35 Medline entries 40 No Medline entry Peptide information for frame 2 45 ORF from 23 bp to 1648 bp; peptide length: 542 Category: similarity to known protein Classification: unclassified 50 I MTLSQAFHLK NNSKKKQMTT EKQKQDANMP KKPVLGSYRG QIVQSKINSF 51 RKPLQVKDES SAATKKLSAT IPKATKPQPV NTSSVTVKSN RSSNMTATTK 101 FVSTTSQNTQ LVRPPIRSHH SNTRDTVKQG ISRTSANVTI RKGPHEKELL 151 QSKTALSSVK TSSSQGIIRN KTLSRSIASE VVARPASLSN DKLMEKSEPV 201 DQRRHTAGKA IVDSRSAQPK ETSEERKARL SEWKAGKGRV LKRPPNSVVT 55

251 QHEPAGQNEK LVGSFWTTMA EEDEQRLFTE KVNNTFSECL NLINEGCPKE 301 DILVTLNDLI KNIPAKKLV KYWICLALIE PITSPIENII AIYEKAILAG 351 AQPIEEMRH IVDILTMKSQ EKANLGENME KSCASKEEVK EVSIEDTGVD

401 VDPEKLEMES KLHRNLLFQD CEKEQDNKTK DPTHDVKTPN TETRTSCLIK 451 YNVSTTPYLQ SVKKKVQFDG TNSAFKELKF LTPVRRSRRL QEKTSKLPDM 501 LKDHYPCVSS LEQLTELGRE TDAFVCRPNA ALCRVYYEAD TT 5 BLASTP hits No BLASTP hits available 10 Alert BLASTP hits for DKFZphtes3_3lalO, frame 2 No Alert BLASTP hits found 15 Pedant information for DKFZphtes3_3lal0, frame 2 Report for DKFZphtes3_3la10.2 20 ELENGTHD 549 61677-36 EMUI [[q] 9.33 [KW] Alpha_Beta 25 EKW3 LOW_COMPLEXITY 2.19 % SEQ DDPQSQHMTLSQAFHLKNNSKKKQMTTEKQKQDANMPKKPVLGSYRGQIVQSKINSFRKP SEG ------30 PRD LQVKDESSAATKKLSATIPKATKPQPVNTSSVTVKSNRSSNMTATTKFVSTTSQNTQLVR SEQ SEG PRD ccccchhhhhhhhhcccccccccceeeeecccccccceeeeec 35 PPIRSHHSNTRDTVKQGISRTSANVTIRKGPHEKELLQSKTALSSVKTSSSQGIIRNKTL SEQ SEG PRD ccccccccccccccccccchhhhhhhhhhccccccccceeeccch SEQ 40 SRSIASEVVARPASLSNDKLMEKSEPVDQRRHTAGKAIVDSRSAQPKETSEERKARLSEW SEG PRD SEQ KAGKGRVLKRPPNSVVTQHEPAGQNEKLVGSFWTTMAEEDEQRLFTEKVNNTFSECLNLI 45 SEG PRD SEQ NEGCPKEDILVTLNDLIKNIPDAKKLVKYWICLALIEPITSPIENIIAIYEKAILAGAQP SEG 50 PRD SEQ IEEMRHTIVDILTMKSQEKANLGENMEKSCASKEEVKEVSIEDTGVDVDPEKLEMESKLH SEG PRD 55 RNLLFQDCEKEQDNKTKDPTHDVKTPNTETRTSCLIKYNVSTTPYLQSVKKKVQFDGTNS SEQ SEG

PCT/IB01/02050

WO 01/98454

	SEQ	AFKELKFLTPVRRSRRLQEKTSKLPDMLKDHYPCVSSLEQLTELGRETDAFVCRPNAALC
	SEG PRD	hhhhhhhhhhhhhhhhhhhcccccccccchhhhhhhhcccc
5	554	
	SEQ	RVYYEADTT
	SEG	• • • • • • • •
	PRD	eeeecccc
10		
	(No	Prosite data available for DKFZphtes3_3lal0-2)
	(No	Pfam data available for DKFZphtes3_3lal0.2)

-375-

DKFZphtes3_31j20

5 group: signal transduction

DKFZphtes3_31j20 encodes a novel 392 amino acid protein that contains a Protein phosphatase 2C motif.

The novel protein shares 95% identity withthe rat protein phosphatase 2C and is expressed ubiquitously. PP2C is a structurally diversified protein phosphatase family with a wide range of functions in cellular signal transduction. The transcription of the PP2Cdelta gene was activated in response to stress, like alcohol or UV irridation. PP2C plays a role in cell cycle control.

The new protein can find application in and the diagnosis/therapy of stress related diseases and cancer, as well as a for 20 modulation of cell cycle and signal transduction.

strong similarity to protein phosphatase 2C (Rattus norvegicus)

25 Sequenced by LMU

Locus: unknown

Insert length: 1436 bp

30 Poly A stretch at pos. 1367, polyadenylation signal at pos. 1341

5

BLAST Results

No BLAST result

10

Medline entries

15 99074314:
Tong Y. Quirion R. Shen SH.; Cloning and characterization of a novel
mammalian PP2C
isozyme. J Biol Chem 1998 Dec 25:273(52):35282-90

20

Peptide information for frame 2

25

55

ORF from 56 bp to 1231 bp; peptide length: 392 Category: strong similarity to known protein Classification: Protein management

30 Prosite motifs: PP2C (147-155)

- I MDLFGDLPEP ERSPRPAAGK EAQKGPLLFD DLPPASSTDS GSGGPLLFDD
 51 LPPASSGDSG SLATSISQMV KTEGKGAKRK TSEEEKNGSE ELVEKKVCKA
 35 101 SSVIFGLKGY VAERKGEREE MQDAHVILND ITEECRPPSS LITRVSYFAV
 151 FDGHGGIRAS KFAAQNLHQN LIRKFPKGDV ISVEKTVKRC LLDTFKHTDE
 201 EFLKQASSQK PAWKDGSTAT CVLAVDNILY IANLGDSRAI LCRYNEESQK
 251 HAALSLSKEH NPTQYEERMR IQKAGGNVRD GRVLGVLEVS RSIGDGQYKR
 301 CGVTSVPDIR RCQLTPNDRF ILLACDGLFK VFTPEEAVNF ILSCLEDEKI
 40 351 QTREGKSAAD ARYEAACNRL ANKAVQRGSA DNVTVMVVRI GH
- BLASTP hits

No BLASTP hits available

Alert BLASTP hits for DKFZphtes3_31j20, frame 2

50 No Alert BLASTP hits found

Pedant information for DKFZphtes3_31j20, frame 2

Report for DKFZphtes3_31j20.2

ELENGTHD 410

44759.85 EMWI [[q] 7.95 [HOMOL] TREMBL: AF095927_1 product: "protein phosphatase 20"; Rattus norvegicus protein phosphatase 20 mRNA, complete cds. EFUNCATI 03.01 cell growth ES. cerevisiae, YDLOObwl be-25 be-25 EFUNCATI D9.16 mitochondrial biogenesis ES. cerevisiae. YDLOObwl be-25 10 IFUNCATI 11.01 stress response ES. cerevisiae, YDLOObwl be-25 EFUNCATD 01.05.04 regulation of carbohydrate utilization cerevisiae, YDLOObwl be-25 EFUNCATD 98 classification not yet clear-cut II. cerevisiae. YERDA9cl le-23 **IFUNCATI** 99 unclassified proteins ES. cerevisiae, YORD9Dcl le-12 ### IFUNCATI 03.22 cell cycle control and mitosis 20 ES. cerevisiae, YJL005w3 3e-10 D3.10 sporulation and germination ES. cerevisiae. **EFUNCATE** YJL005wJ 3e-10 **EFUNCATE** 30.02 organization of plasma membrane ES. cerevisiae, YJL005w] 3e-10 EFUNCATI 01.03.10 metabolism of cyclic and unusual nucleotides ES. cerevisiae, YJLOO5wl 3e-lo **EFUNCATI** 10.04.03 second messenger formation ES. cerevisiae, YJL005w3 3e-10 PR01023F 30 [BLOCK2] PROOL77D **EBFOCK23** BF070351 [BFOCK2] [Brockz] BF07035H BL070356 **EBFOCK21** EBLOCKSI BLO1032C Protein phosphatase 2C proteins
EBLOCKSI BLO1032B Protein phosphatase 2C proteins 35 dlabq___ 4.98.1.1.1 Protein serine/threonine **EZCOPI** phosphatase 2C [Huma le-107 3.1.3.43 [Pyruvate dehydrogenase (lipoamide)]-EECI phosphatase 3e-09 40 3.1.3.16 Phosphoprotein phosphatase 7e-35 EECI 4.6.1.1 Adenylate cyclase 2e-11 EECI duplication 5e-ll [PIRKW] **EPIRKWI** tandem repeat &e-09 45 serine/threonine-specific phosphatase 2e-27 **EPIRKU** magnesium be-26 **EPIRKW3 EPIRKWI** camp biosynthesis 5e-11 **CPIRKWI** liver 2e-27 **EPIRKWI** leucine zipper le-Ob 50 mitochondrion 3e-09 [PIRKU] phosphoric monoester hydrolase 7e-35 **EPIRKWI EPIRKU**I phosphorus-oxygen lyase 2e-11 ESUPFAMI leucine-rich alpha-2-glycoprotein repeat homology 2e-11 55 ESUPFAMI yeast adenylate cyclase catalytic domain homology 2e-ll

kinase interaction domain homology 3e-11

yeast adenylate cyclase 5e-11

ESUPFAMI ESUPFAMI

CPROSITED PP2C 1 Protein phosphatase 2C [PFAM] [KW] Alpha_Beta 5 AARGLSVCRCCRLHPASAMDLFGDLPEPERSPRPAAGKEAQKGPLLFDDLPPASSTDSGS SEQ PRD GGPLLFDDLPPASSGDSGSLATSISQMVKTEGKGAKRKTSEEEKNGSEELVEKKVCKASS SEQ 10 PRD VIFGLKGYVAERKGEREEM@DAHVILNDITEECRPPSSLITRVSYFAVFDGHGGIRASKF SEQ PRD AAQNLHQNLIRKFPKGDVISVEKTVKRCLLDTFKHTDEEFLKQASSQKPAWKDGSTATCV 15 SEQ PRD LAVDNILYIANLGDSRAILCRYNEESQKHAALSLSKEHNPTQYEERMRIQKAGGNVRDGR. SEQ eeccceeeeccccceeeeeeccccccceeeee PRD 20 VLGVLEVSRSIGDGQYKRCGVTSVPDIRRCQLTPNDRFILLACDGLFKVFTPEEAVNFIL SEQ PRD SCLEDEKIRTREGKSAADARYEAACNRLANKAVQRGSADNVTVMVVRIGH SEQ 25 PRD Prosite for DKFZphtes3_31j20.2 30 PD0C00792 165->174 **PP2C** 5E01035 35 ----HMM_NAME Protein phosphatase 20 40 HMM *G1CcMQGPRWRMsMEDaHiaylNF....pcnlDWWhiMFFGVFDGHg +++ +G R++M+DAH+ + ++ +++F+VFDGHG YVAERKG--EREEMQDAHVILNDITEECRPPSSLITR-758 Query VSYFAVFDGHG 45 173

189

GDQCSQWCgeHWHdII*

G+++S++ ++++++ +

174 GIRASKFAAQNLHQNL

HMM

50

Query

WO 01/98454 DKFZphtes3_5k22

5 group: signal transduction

DKFZphtes3_5k22 encodes a novel 455 amino acid protein with similarity to human paraneoplastic neuronal antigen MAl-

- 10 Antibodies against MAL where found in patients with paraneoplastic neurological disorders. The protein is predominantly expressed in testis and brain, but ESTs are also found in liver, lung uterus and kidney.
- 15 The new protein can find application in studying/therapy of paraneoplastic neurological disorders.

strong similarity to paraneoplastic neuronal antigen MAI

20 Sequenced by @iagen

Locus: unknown

25 Insert length: 3534 bp
Poly A stretch at pos. 3514, polyadenylation signal at pos. 3494

PCT/IB01/02050 WO 01/98454 1501 GTACAGTGCA TCAACCCCTC CAACCTGCTC TTGGCCAAGG AGACAAAAGA
1551 GATATTGGAA GGAGGGGAAA GAGAAGCCCA GACAAACAGC AGATGAGTTG
1601 AGTGGGGCAG AGGGACAGGG CAGCCAGACC AAGGCCAAGC CTTCTCACCC
1651 TTGGCCAGCT GGAAGGGACT TCAGCAACCA AGACCACCTG GCAACAGGCT
1701 CAGTGGGGGT CAGGTCCAGG TCCCCGAAGA GGTGCTGGAG AGGAAAGCAG
1751 GGAGCCACTG CATCCAGCAC ATGGGGTGCC TGGGCCTCAG ATGGGGACCC
1801 CAAAGAAGCA GAAGCTGAAG AAGGTACGG TGGGGGTTCT GTCCTGCTCA
1851 TCCAACCACC CCTAAATACC CACCCTGTGG ACTTTGAGCT GAACATGCCC 5 1901 ACTGGCCCCC AGGCCACATG GGACCTGGAG GAGCCTACCT GGGGCCTGCC 1951 CCTGCCAGCA GGTGCCAGGG CTGGTGAGGA AGAGCTGGGG GGCAGAGGTA 10 2001 AAGCCCTGCA GGGGAGGCCA CAGGGTCCAT CCCGTCTTCA GGATCATCTA 2051 CACTGCACTA GGGGAGCCCC AGGAAGGCAG CACCCTGGAG GCCCTGTGCC ZIDI AGTGAGGACA GGAGACCCTA AGGCCCCGGG AGCCCAGTGC CAGCCAGAGG 2151 TTGTGCAGGC AAGGAGACCA AAGATTGATG AGAAGACCCC CAGCAGGGGT 2201 ACTGGGTACC CGGCAGGCCA GTGCCCTCAC AGTTGACTTG GACCAGGGTG 15 2251 GCTGTGAAGG GAAGTCTTTG TTGCAAAGGA GGAGGAGGAA AAGGGAGGAC TOAGACOACO GODACATATT TTTGTTTCTT CTGTTTT CTGDATGGTT TOAGACT 2351 CCTGGAGAGA TCAAGCAAGG AGAACCTGGG GCTGCCATGG CCAAAGCAAC 2401 TCAACAGATG CCAATGCCAA TTCCAAGGCC AGCCACAACC CTGCCACCTT 2451 GGGGAATCCA GCCTGGAGGC ATCCCCTAAG CAGCCAGCCA TGGCCTGGGT 20 25D1 GGAGGCACCT GAAGACGTCT GTCCCAAACT CCCCCAGCCC TGAGCTGGGA 2551 GATGACAGGG GGAAAGAGGC CCTCTCAAGG GTGCCAGATG CCTGGGTCTC 2601 CCAAGAGGGG TCCCCCAACT CACCGTTCCC GGGACAGGCT GCCCCCTGTT 2651 CCAGGAAGCT CATCCTCACC TGTGTAGGCC CCTGTAGTGA CCCACGCGTC 2701 CAGCAGACGC CCACCCACCG CTAGCCGTTG TTCCTGTGCA AAGTAGTGTG 25 2751 CTATGCACCC ACCCAGGTGG CCGCCTCTGG GCCCAAGGCA CATGCTGTGA 2801 GCTTCCTGTG AGCCCAGGCT CTGCTCACTG CTGTCCCGCG TCATGAGCAC 2851 CACCTCTGCT TTCCCTGTGT AGATCTAGGC CAGTGGCTGC TTGTTCTTGT 2901 GGAGCTGTGT GTGTTCTTCT CTGAGCAGCT CCTCCCCGA GTCCCCCAGC 2951 ACAGTCCCAG GAGATGACAG GAAGGAAGCA CCAGGGCAAG GCGGACGCTC 30 **3001 ACCCTGTGAC CACGATGGTG ACCGTGGCTG TGGGAGGAAG AACTGGACCC** AGGACGAGC CACGGCTGCCC TGCCTGAGGC TCCCGAGGAG CTTTGTGCTT

ATGGTTTCA CCCCTGTTGT TACTCATAC TCATTGTGGT TCCCTGTGT TCCCTGTG 3201 CATAGGGCAG GGCCCTGCCC- CAGCAGATGG -GCTTGGGAGG GGGCTCCCTA 35 3251 AAGCCAGTGG ACACTGCCAG AGTCTACCTT CCTGGCAAGA GGCAGACCCC 40 3501 CTTCATTCAG TGTTAAAAAA AAAAAAAAAA AAAA

BLAST Results

45

No BLAST result

50

Medline entries

99158179:

Mal, a novel neuron- and testis-specific protein, is recognized by the serum of patients with paraneoplastic neurological disorders.

Peptide information for frame 1

ORF from 229 bp to 1593 bp; peptide length: 455 Category: strong similarity to known protein Classification: unclassified

10 I MPLTLLQDUC RGEHLNTRRC MLILGIPEDC GEDEFEETLQ EACRHLGRYR
51 VIGRMFRREE NAQAILLELA QDIDYALLPR EIPGKGGPWE VIVKPRNSDG
101 EFLNRLNRFL EEERRTVSDM NRVLGSDTNC SAPRVTISPE FUTWAQTLGA
151 AVQPLLEQML YRELRVFSGN TISIPGALAF DAWLEHTTEM LQMWQVPEGE
201 KRRRLMECLR GPALQVVSGL RASNASITVE ECLAALQQVF GPVESHKIAQ
15 251 VKLCKAYQEA GEKVSSFVLR LEPLLQRAVE NNVVSRRNVN QTRLKRVLSG
301 ATLPDKLRDK LKLMKQRRKP PGFLALVKLL REEEEWEATL GPDRESLEGL
351 EVAPRPPARI TGVGAVPLPA SGNSFDARPS QGYRRRRGRG QHRRGGVARA
401 GSRGSRKRKR HTFCYSCGED GHIRVQCINP SNLLLAKETK EILEGGEREA

451 QTNSR

5

20

BLASTP hits

25 No BLASTP hits available

Alert BLASTP hits for DKFZphtes3_5k22, frame 1

TREMBLNEW:ABO20690_1 gene: "KIAAO883"; product: "KIAAO883

30 protein";

Homo sapiens mRNA for KIAAO883 protein; complete cds.; N = 1;

Score =

722; P = 2.4e-71

- 35 TREMBL:AFD37364_1 gene: "MAl"; product: "paraneoplastic neuronal antigen MAL"; Homo sapiens paraneoplastic neuronal antigen MAL (MAL)

 mRNA, complete cds., N = 1, Score = 665, P = 2.6e-65
- 40
 >TREMBLNEW:ABD20690_1 gene: "KIAAD883"; product: "KIAAD883
 protein"; Homo
 sapiens mRNA for KIAAD883 protein; complete cdsLength = 364

45 HSPs:

Score = 722 (108.3 bits), Expect = 2.4e-71, P = 2.4e-71 Identities = 156/348 (44%), Positives = 215/348 (61%)

Query: 1
MPLTLLQDWCRGEHLNTRRCMLILGIPEDCGEDEFEETLQEACRHLGRYRVIGRMFRREE 60
M L LL+DWCR ++ ++ ++ GIP D E E +E LQE +
LGRYR++G++FR++E

55 Shict: 1

55 Sbjct: 1 MALALLEDWCRIMSVDEQKSLMVTGIPADFEEAEIQEVLQETLKSLGRYRLLGKIFRKQE 60

Query: 61
NAQAILLELAQDIDYALLPREIPGKGGPWEVIVKPRNSDGXXXXXXXXXXXXXXXXXXXXDDM 120

NAQAILLELAQDIDYALLPREIPGKGGPWEVIVKPRNSDGXXXXXXXXXXXXXXXXXTVSDM 13C

M ZVT

5 Sbjct: 61
NANAVLLELLEDTDVSAIPSEVQGKGGVWKVIFKTPNQDTEFLERLNLFLEKEGQTVSGM 120

Query: 121 NRVLGSDTNCSAPRVTISPEFWTW--AQTLGAAVQPLLEQMLYRELRVFSGNTISIPGAL 178

10 R LG + A ISPE Q + A QPLL M YR+LRVFSG+ + P Sbjct: 121 FRALGQEGVSPATVPCISPELLAHLLGQAMAHAPQPLLP-

MRYRKLRVFSGSAVPAPEE 379

15 Query: 179
AFDAWLEHTTEMLQMWQVPEGEKRRRLMECLRGPALQVVSGLRASNASITVEECLAALQQ 238
+F+ WLE TE+++ W V E EK+R L E LRGPAL ++ ++ A N
SI+VEECL A +Q

Sbjct: 180
SFEVWLEQATEIVKEWPVTEAEKKRWLAESLRGPALDLMHIVQADNPSISVEECLEAFKQ 239

25 L++V+
Sbjct: 240
VFGSLESRRTAQVRYLKTYQEEGEKVSAYVLRLETLLRRAVEKRAIPRRIADQVRLEQVM 299

Query: 299 SGATLPDKLRDKLKLMKQRRKPPGFLALVKLLREEEEWEATLGPDRESLE 30 348

+GATL L +L+ +K + PP FL L+K++REEEE EA+ + ES+E Sbjct: 300 AGATLNQMLWCRLRELKDQGPPPSFLELMKVIREEEEEEASF--ENESIE 347

- -

35

20

Pedant information for DKFZphtes3_5k22, frame 1

Report for DKFZphtes3_5k22.1

40

ELENGTHD 455
EMWD 51514.34
Epil 9.27

45 EHOMOLI TREMBLNEW: ABD2Db9D_1 gene: "KIAAD&83"; product: "KIAAD&83 protein"; Homo sapiens mRNA for KIAAD&83 protein; complete cds. 3e-75
EBLOCKSI BLDB87bB Indoleamine 2-3-dioxygenase proteins
EPFAMI Zinc finger; CCHC class

50 [KW] Alpha_Beta [KW] LOW_COMPLEXITY 13.41 %

SEQ NAQAILLELAQDIDYALLPREIPGKGGPWEVIVKPRNSDGEFLNRLNRFLEEERRTVSDM

WO 01/98454 PCT/IB01/02050 SEG PRD SEQ NRVLGSDTNCSAPRVTISPEFWTWAQTLGAAVQPLLEQMLYRELRVFSGNTISIPGALAF 5 SEG PRD SEQ DAWLEHTTEMLQMWQVPEGEKRRRLMECLRGPALQVVSGLRASNASITVEECLAALQQVF SEG 10 PRD SEQ GPVESHKIAQVKLCKAYQEAGEKVSSFVLRLEPLLQRAVENNVVSRRNVNQTRLKRVLSG SEG -----xxxxxxxxxxxxxxx..... PRD 15 SEQ ATLPDKLRDKLKLMKQRRKPPGFLALVKLLREEEEWEATLGPDRESLEGLEVAPRPPARI SEG PRD 20 SEQ TGVGAVPLPASGNSFDARPSQGYRRRRGRGQHRRGGVARAGSRGSRKRKRHTFCYSCGED SEG ····· PRD SEQ GHIRVQCINPSNLLLAKETKEILEGGEREAQTNSR 25 SEG ceeeeecccchhhhhhhhhhhhcccccccccc (No Prosite data available for DKFZphtes3_5k22.1) 30 Pfam for DKFZphtes3_5k22-1 35 HMM_NAME - Zinc finger - CCHC - class --

> *@kCWNCGKPGHMMRDCPE* C++CG+ GH+ +C +

TFCYSCGEDGHIRVQCIN

429

412

HMM

40

Querv

DKFZphtes3_7nl2

5 group: transmembrane protein

DKFZphtes3_7nl2 encodes a novel 703 amino acid protein without similarity to known proteins.

- 10 The novel protein contains 1 transmembrane domain
 No informative BLAST results; No predictive prosite, pfam or SCOP
 motife.
- The new protein can find application in studying the expression profile of testis-specific genes and as a new marker for testicular cells.

putative protein

20

contains transmembrane domain perhaps complete cds.

Sequenced by BMFZ

25

30

Locus: unknown

Insert length: 2347 bp
Poly A stretch at pos. 2271, polyadenylation signal at pos. 2253

1 CGGCTGCAGT CTGGGCCGGG GCCCTGTGCC GCTGAAGACA TGGAGTTTGT 51 GTCTGGATAC CGGGATGAGT TCCTTGATTT CACTGCCCTT CTCTTCGGCT LOL GGTTCCGAAA GTTTGTGGCA GAGCGTGGAG CTGTAGGGAC TAGCCTTGAG 35 151 GGCCGCTGCC GGCAGCTGGA GGCCCAGATC AGAAGGCTAC CCCAGGACCC **2DL TGCCCTTTGG GTGCTCCATG TCCTGCCCAA CCATAGTGTG GGCATCAGCC** 251 TGGGGCAAGG GGCAGAACCA GGTCCTGGAC CAGGCCTGGG GACTGCCTGG BD1 CTCCTGGGAG ACAACCCTCC ACTCCACCTG CGAGACCTGA GCCCCTACAT 351 CAGCTTTGTC AGCCTAGAGG ATGGGGAGGA AGGGGAGGAG GAAGAGGAGG 40 4D1 AAGATGAAGA AGAAGAGAAG AGAGAGGACG GGGGTGCAGG CAGCACAGAG 451 AAGGTGGAAC CAGAGGAGGA CCGGGAGCTA GCCCCTACCA GCAGGGAGTC 501 CCCCCAGGAA ACAAACCCTC CAGGAGAGTC AGAGGAGGCT GCCCGGGAGG 551 CAGGAGGTGG CAAGGATGGC TGCCGAGAGG ACAGGGTGGA GAACGAAACA **LOS AGACCCCAGA AGAGGAAGGG ACAGAGGAGT GAGGCTGCCC CCCTGCACGT** 45 **L51 TTCCTGTCTC TTACTTGTGA CGGATGAGCA TGGCACCATC TTGGGCATTG** 7D1 ATCTGCTAGT GGATGGAGCC CAGGGAACCG CAAGCTGGGG CTCAGGGACC 75% AAGGACCTGG CTCCTTGGGC CTATGCTCTC CTCTGTCACA GCATGGCCTG BD1 TCCCATGGGC TCTGGGGATC CCCGAAAGCC CCGACAGCTT ACTGTGGGAG 851 ATGCCCGGCT GCATCGAGAG CTGGAGAGCT TGGTCCCAAG GCTAGGTGTG 9D1 AAGTTAGCCA AAACCCCAAT GCGGACATGG GGTCCCCGGC CAGGCTTCAC 50 951 CTTTGCTTCC CTTCGTGCTC GAACCTGCCA TGTGTGTCAC AGGCACAGCT 1001 TTGAAGCGAA GCTGACACCT TGCCCCCAGT GTAGTGCTGT CTTGTATTGT 1051 GGAGAGGCTT GTCTCCGGGC TGACTGGCAG CGGTGCCCAG ATGATGTGAG INDI TCACCGATTT TGGTGCCCAA GGCTTGCAGC CTTCATGGAG CGGGCAGGAG 1151 AACTGGCAAC CCTACCTTTT ACCTACACCG CAGAGGTGAC CAGTGAAACC
1201 TTCAACAAAG AGGCCTTCCT GGCCTCTCGG GGCCTCACTC GTGGCTATTG
1251 GACCCAGCTC AGCATGCTGA TTCCAGGCCC GGGCTTCTCC AGACACCCCC 55 LBD1 GAGGCAACAC GCCATCCCTC AGCCTTCTTC GCGGTGGAGA CCCCTACCAG

	WO 01/98454				PCT/IB01/02050	
	1351 CTTCTCCAGG 1401 ACCCCGGGGT 1451 TGAAGATCCA 1501 TTTTGGGAGC	GTTTTTGTCC CGTGGTGGAG	CTGAGCTCAA	CATCCAAAAC AGTTTGACCT	CCCCACATCC AAACAGTCAC TGTCATGGTG AGCTGCAGTT	
5	1551 TGTAGGTGAT 1601 AGAGGGACAG 1651 CGGCCCAGCT	GGCCTGCCCC CCTGGAGGTG CTGGCACTAA	CCGAAAGCGA	CGAGCAGCAT CTGGTTCCGG GGCCGCAGGG	TTTACCCTGC CATATCAGCA ACCTGCAGAT AAGCCTGACC	
10	1751 TGGTTATTGG 1801 TCTCTGCCCC 1851 CAGCGAGTAC 1901 GAGGGGGCAC 1951 CTCAGAGCGG	ATTTAACTCC GGTTACAGTC AGCTGTGTGA CAGCCCTCCC	GGGTTTGCTC CCTCCGAGTG TGGACGGCCA CAGCCCAACC	TCAAGGATAC CCAGCCTTCT GACCATGGCG CCTTCCGCTC	GTGGCTGAGG TCACCGAGAG GTGGCCACTG CCCCTTTCGC	
15	2001 CCACCTGGTT 2051 GGCCCCCACC 2101 AGGCGCCGAG 2151 AATGCTGATA 2201 TGAAAACACT	CCCATCCCA GAGAAAAGAA CCCTAGTAGT CAAGGCCTAG	ACTCCTCTG ACCTGGGCGG CCCCAGCTCC GGGGAGGACA	CTCCTCCTGC GGGGCCCGCC CAAACACTGA GGTTGGTAAA	CCCACCGA GGCGGAAATG AAGGAAAACG ACATGAAAAG	
20	AAAAAAAA 2255 AAAAAAAA 40E5					
	BLAST Results					
No BLAST result						
30		Medline entries				
50	No Medline entry					
	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,					
Peptide information for frame 1						
ORF from 4D bp to 2148 bp; peptide length: 703 40 Category: putative protein Classification: Transmembrane proteins unclassified						
45		VLPNHSVGIS DGEEGEEEEE PGESEEAARE	LGQGAEPGPG EDEEEEKRED AGGGKDGCRE	PGLGTAWLLG	DNPPLHLRDL PEEDRELAPT KRKGQRSEAA	
50	351 DDVSHRFWCP 401 RGYWTQLSML : 451 PPHPPRGVFV	RTCHVCHRHS RLAAFMERAG IPGPGFSRHP PELNIQNKQS	FEAKLTPCP@ ELATLPFTYT RGNTPSLSLL LKIHVVEAGK	CSAVLYCGEA AEVTSETFNK RGGDPYQLLQ EFDLVMVFWE	CLRADW@RCP EAFLASRGLT GDGTALMPPV LLVLLPHVAL	
55	503 ELQFVGDGLP (553 DLQIKVSARP (603 FTESSEYSCV (653 AFIFHLVYKP (YHLFQGPKPD	LVIGFNSGFA GGGTSPPQPN	LKDTWLRSLP PFRSPFRLRA	RLQSLRVPAF ADNCMSWYCN	

BLASTP hits

No BLASTP hits available Alert BLASTP hits for DKFZphtes3_7nl2, frame 1 No Alert BLASTP hits found **10** . Pedant information for DKFZphtes3_?nl2, frame 1 Report for DKFZphtes3_7nl2.1 15 703 ELENGTHI EWWI 77375.75 [[q] 6.45 TRANSMEMBRANE 20 [KW] [KW] LOW_COMPLEXITY 15.22 % SEQ MEFVSGYRDEFLDFTALLFGWFRKFVAERGAVGTSLEGRCRQLEAQIRRLPQDPALWVLH 25 SEG PRD MEM **VLPNHSVGISLGQGAEPGPGPGLGTAWLLGDNPPLHLRDLSPYISFVSLEDGEEGEEEE** SEQ 30 SEG --------XXXXXXXX PRD MEM **EDEEEEKREDGGAGSTEKVEPEEDRELAPTSRESPQETNPPGESEEAAREAGGGKDGCRE** SEQ 35 SEG _xxxxxxxxxxx...... PRD MEM DRVENETRP&KRKG&RSEAAPLHVSCLLLVTDEHGTILGIDLLVDGA>ASWGSGTKDL SEQ 40 SEG PRD MEM SEQ APWAYALLCHSMACPMGSGDPRKPRQLTVGDARLHRELESLVPRLGVKLAKTPMRTWGPR 45 SEG PRD MEM SEQ PGFTFASLRARTCHVCHRHSFEAKLTPCPQCSAVLYCGEACLRADWQRCPDDVSHRFWCP 50 SEG PRD MEM SEQ RLAAFMERAGELATLPFTYTAEVTSETFNKEAFLASRGLTRGYWTQLSMLIPGPGFSRHP 55 SEG PRD MEM

WO 01/98454 PCT/IB01/02050 SEQ RGNTPSLSLLRGGDPYQLLQGDGTALMPPVPPHPPRGVFVPELNIQNKQSLKIHVVEAGK SEG PRD CCCCCCEEECCCCCCCCCCCCCCCCCCCCCCChhhhhheeeeeccc MEM 5 SEQ EFDLVMVFWELLVLLPHVALELQFVGDGLPPESDEQHFTLQRDSLEVSVRPGSGISARPS SEG PRD ---mmmmmmmmmmmmm...... MEM 10 SGTKEKGGRRDL@IKVSARPYHLF@GPKPDLVIGFNSGFALKDTWLRSLPRL@SLRVPAF SEQ SEG PRD MEM 15 SEQ FTESSEYSCVMDGQTMAVATGGGTSPPQPNPFRSPFRLRAADNCMSWYCNAFIFHLVYKP SEG ----x PRD MEM 20 AQGSGARPAPGPPPPSPTPSAPPAPTRRRRGEKKPGRGARRRK SEQ SEG PRD cccccccccccccccchhhhhccccccccccc MEM 25 (No Prosite data available for DKFZphtes3_7nl2.1) (No Pfam data available for DKFZphtes3_7nl2.1) 30 DKFZphtes3_9el6 35 group: transmembrane protein DKFZphtes3_9el6 encodes a novel 539 amino acid protein without similarity to known proteins. The novel protein contains 1 transmembrane region. The only EST 40 described so far is from testis. No informative BLAST results: No predictive prosite, pfam or SCOP motife. 45 The new protein can find application in studying the expression profile of testis-specific genes and as a new marker for testicular cells. 50 putative protein 1 EST hit perhaps complete cds. 55 Sequenced by DKFZ Locus: unknown

Insert length: 2011 bp

Poly A stretch at pos. 1986, no polyadenylation signal found

```
5
                             L CATGGCAACA TGAGCAGTGC TGAGATAATT GGTTCTACAA ATCTTATAAT
                          51 TCTGCTAGAG GATGAAGTCT TTGCCGATTT TTTCAACACA TTTCTTTCCC
                       LOL TCCCGGTTTT TGGTCAGACA CCATTTTATA CTGTTGAAAA TTCACAGTGG
                    BSL AGCTTGTGGC CAGAAATACC TTGTAACTTG ATTGCCAAAT ACAAAGGGTT 20L ATTGACCTGG TTGGAAAAAT GCCGATTACC TTTCTTCTGT AAAACAAACT 25L TGTGTTTCCA TTACATTCTC TGTCAGGAGT TCATCAGTTT CATTAAGTCC 30L CCAGAAGGAG CCAAGATGAT GAGATGGAAA AAGGCAGACC AGTGGCTACT 35L CCAGAAATGC ATTGGCGGGG TCAGAGGGAT GTGGCGCTTC TATTCCTACC 40L TCACAGGCAG TGCAGGTGAA GAATTGGTGG ATTTCTGGAT CCTTGCTGAG 45L AACATCCTGA GCATAGATGA GATGGACCTG GAAGTGAGA GCTCCAGGG 55L TGGTAACCTT CTGTAACATG AACATCAAGT CCCTCCTGAA CCTCTCCATC 60L TGGCATCCCA ACCAATCAAC CACTAGGAGG GAGATCCTGA GCCACATGCA 65L GAAAGTGGCT CTGTTCAAAC TCCAGAGGTA TTGGCTTCCC AACTTTTACA 70L CCCACACCAA GATGACCATG GCCAAGGAGG AAGCATGCCA TGGTCTGATG 75L CAAGAGTACG AGACTCGCTT ATACAGCGT TGCTACACCC ACATAGGAGG 80L GCTCCCTCTG AACATCGAAC TCCAGAGAGG AAGCATGCCA TGGTCTGATG 75L CAAGAGTACG AGACTCGCTT ATACAGCGT TGCTACACCC ACATAGGAGG 80L GCTCCCTTG AACATGGACA TCAAGAAGTG CCACCACTTT CAGAAACGGT 85L ACTCCAGCAG GAAAGCCAAG AGGAAGATGT GGCAATTGGT AGATCCTGAC 90L TCTTGGTCTC TGGAAATGGA TCTCAAGCCA GATGCTATTG GTATGCCCCT 95L ACAGGAGACA TGTCCTCAAG AGAAGGTGGT TATACAAATG CCTTCCCTGA ACATGGACA ACAAGGAATCA GTTCCCCTGAA ACAAGGAATCA GTTCCCCTGGA AAAGGGATATG
                       151 AGCTTGTGGC CAGAAATACC TTGTAACTTG ATTGCCAAAT ACAAAGGGTT
  10
 15
 20
                P51 ACAGGAGACA TGTCCTCAAG AGAAGGTGGT TATACAAATG CCTTCCCTGA
D01 AAATGGCTTC TTCAAAGGAA ACAAGAATCA GTTCCCTGGA AAAGGATATG
D051 CATTATGCAA AAATATCCAG CATGGAGAAT AAAGCCAAGA GCCACCTCCA
L101 CATGGAAGCC CCCTTTGAGA CAAAGGTCTC TACCCACCTG AGGACTGTCA
L201 CAAAGCTTCT CAATCACTCC TCCAAGATGA CAATTCAGAA GGCCATCAAG
L201 CAAAGCTTCT CCTTAGGATA CATCCACTTG GCCTTGTGG CTGATGCCTG
L251 TGCAGGGAAC CCTTTCCGGG ACCACCTGAA GAAGCTGAAT TTGAAAGTGG
L301 AGATCCAACT TCTTGACCTC TGGCAGGACT TGCAGCATTT CCTCAGTGTC
L351 CTTCTGAATA ACAAAAAGAA TGGGAATGCA ATCTTTCGTC ACTTGCTGGG
L401 TGACAGAATC TGCGAGCTCT ACCTGAATGA GCAGATTGGT CCGTGCTTAC
L451 CACTCAAATC CCAAACCATT CAGGGCCTGA AGGAACTATT GCCCTCTGGG
L501 GATGTGATCC CCTGGATTCC CAAAGCCCAG AAGGAGATTT GCAAGATGCT
L551 CAGTCCCTGG TATGATGAGT TTCTAGATGA AGAGGACTAC TGGTTTCTCC
 25
 30
 35
                 1551 CAGTCCCTGG TATGATGAGT TTCTAGATGA AGAGGACTAC TGGTTTCTCC
1601 TTTTTACGGT AGGAAGGACT TTGGGTTAGG AAGGAATCAT GAGGATGAGG
                 1651 GAAGAAGAAA GAGTAATTAC TGTTTTAAAA GGGTTATGTG TTAAAGTAAA
                 1701 TGAAATTGTT ATTTTTCCTA GAGTCAACCA AAGATCAGCA TGGTCCCTGT
40
                 1751 TGTTCTAAAG CTAAACCTCT CAAGGAAAAG GACTCAGTGC ATAAGATGAC
                 18D1 TTTGGTGAAA CCCCGTCTCT ACTAAAAATA CAAAAAATTA GCCGGGCGTA
                 1851 GTGGCGGGCG CCTGTAGTCC CAGCTACTTG GGAGGCTGAG GCAGGAGAAT
                 1901 GGTGTGAACC CGGGAGGCGG AGCTTGCAGT GAGCCGAGAT CCCGCCACTG
                 1951 CACGCCAGCC TGGGCGACAG AGCGAGACTC CGTCTCAAAA AAAAAAAAA
45
          AAAAAAAAA G
```

BLAST Results

50

No BLAST result

55

Medline entries

No Medline entry

Peptide information for frame 1

ORF from 10 bp to 1626 bp; peptide length: 539 Category: putative protein Classification: no clue

- 10 LEMDLKPDAI GMPLQETCPQ EKVVIQMPSL KMASSKETRI SSLEKDMHYA
 35 KISSMENKAK SHLHMEAPFE TKVSTHLRTV IPIVNHSSKM TIQKAIKQSF
 401 SLGYIHLALC ADACAGNPFR DHLKKLNLKV EIQLLDLWQD LQHFLSVLN
 451 NKKNGNAIFR HLLGDRICEL YLNEQIGPCL PLKSQTIQGL KELLPSGDVI
 202 ON CONTROL CONTRO
 - BLASTP hits

25 No BLASTP hits available

Alert BLASTP hits for DKFZphtes3_9elb, frame l

30. No Alert BLASTP hits found

Pedant information for DKFZphtes3_9elb, frame l

Report for DKFZphtes3_9elb-1

ELENGTHD 542

40 - 40PS4 EWM3

40 EpII 8-35

EKW3 Alpha_Beta

- - SEQ IAKYKGLLTWLEKCRLPFFCKTNLCFHYILCQEFISFIKSPEGAKMMRWKKADQWLLQKC
- 50 SEQ IGGVRGMWRFYSYLTGSAGEELVDFWILAENILSIDEMDLEVRDYYLSLLLMLRATHLQE PRD ccccceeeeeecccccchhhhhhhhhhhhccc
 - SEQ GSRVVTLCNMNIKSLLNLSIWHPNQSTTRREILSHMQKVALFKLQSYWLPNFYTHTKMTM

SEG ZMZFEWDFKbDaiewbfgelcbgeknnigwb2rwazzkelbizzfekdwhakizzwen

- PRD ccccccccccccccccceeeeecccccccccchhhhhh
- SEQ KAKSHLHMEAPFETKVSTHLRTVIPIVNHSSKMTIQKAIKQSFSLGYIHLALCADACAGN
- - SEQ PFRDHLKKLNLKVEIQLLDLWQDLQHFLSVLLNNKKNGNAIFRHLLGDRICELYLNEQIG
- 10 SEQ PCLPLKSQTIQGLKELLPSGDVIPWIPKAQKEICKMLSPWYDEFLDEEDYWFLLFTVGRT
 - PRD cccccchhhhhhhhccccccceeeccchhhhhhhcccchhhhhccccceeeccccc

SEQ LG

PRD cc

15

(No Prosite data available for DKFZphtes3_9elb.1)

(No Pfam data available for DKFZphtes3_9elb.1)

20

25

The PROSITE is a database of protein families and domains. It consists of biologically significant sites, patterns and profiles that help to reliably identify to which known protein family (if any) a new sequence belongs. World Wide Web URL http://www.expasy.ch/prosite/ is the entry point to the database. A description of the prosite consensus patterns follows.

30

NAME: N-glycosylation site.
CONSENSUS: N-{P}-ESTI-{P}.

NAME: Glycosaminoglycan attachment site.

35 CONSENSUS: S-G-x-G-

NAME: Tyrosine sulfation site.

NAME: cAMP- and cGMP-dependent protein kinase

40 phosphorylation site.

CONSENSUS: ERKI(2)-x-ESTI.

NAME: Protein kinase C phosphorylation site.

CONSENSUS: EST3-x-ERK3.

45

NAME: Casein kinase II phosphorylation site-

CONZENZUS: EZTJ-x(2)-EDEJ-

NAME: Tyrosine kinase phosphorylation site.

NAME: N-myristoylation site.

CONSENSUS: G-{EDRKHPFYW}-x(2)-ESTAGCNJ-{P}.

55 NAME: Amidation site.
CONSENSUS: x-G-ERKI-ERKI.

NAME: Aspartic acid and asparagine hydroxylation site.

. CONSENSUS: $C-x-\mathbb{E}DNJ-x(4)-\mathbb{E}FYJ-x-C-x-C$.

NAME: Vitamin K-dependent carboxylation domain.

CONZENZUZ: $\times (J2)-E-\times (3)-E-\times-C-\times (b)-\mathbb{C}DENJ-\times-\mathbb{E}LIVMFYJ-\times (9)-$

5 EFYW3.

NAME: Phosphopantetheine attachment site.

CONSENSUS: EDEQGSTALMKRHJ-ELIVMFYSTACJ-EGNQJ-ELIVMFYAGJ-

EDNEKHZJ-Z-ELIVMZTJ-

10 CONSENSUS: {PCFY}-ESTAGCPQLIVMFJ-ELIVMATNJ-EDENQGTAKRHLMJ-

ELIVMUSTAI-ELIVGSTACRI-

CONSENSUS: x(2)-[LIVMFA].

NAME: Acyl carrier protein phosphopantetheine domain

15 profile.

NAME: Prokaryotic membrane lipoprotein lipid attachment

site.

CONSENSUS: {DERK}(b)-ELIVMFWSTAGI(2)-ELIVMFYSTAGCQI-EAGSI-C.

20

NAME: Prokaryotic N-terminal methylation site.

CONSENSUS: EKAHEQSTAGD-G-ELAVIJATO-ESTD-ELTD-ELIVPD-E-

ELIVMFWSTAGI(14).

25 NAME: Prenyl group binding site (CAAX box).

CONSENSUS: C-{DENQ}-ELIVMI-x>.

NAME: Protein splicing signature.

COTED - CTEANULE - CLIVALE - CLIVALE

30

NAME: Endoplasmic reticulum targeting sequence.

CONSENSUS: EKRHQSAJ-EDENQJ-E-L>.

NAME: Microbodies C-terminal targeting signal.

35 CONSENSUS: ESTAGONI-ERKHI-ELIVMAFYI>.

NAME: Gram-positive cocci surface proteins 'anchoring'

hexapeptide.

CONSENSUS: L-P-x-T-G-ESTGAVDED.

40

NAME: Bipartite nuclear targeting sequence.

NAME: Cell attachment sequence.

CONSENSUS: R-G-D.

45

NAME: ATP/GTP-binding site motif A (P-loop).

CONZENZUZ: [AG]-x(4)-G-K-[ZT].

NAME: Cyclic nucleotide-binding domain signature 1.

50 CONSENSUS: ELIVMI-EVICD-x(2)-G-EDENGTAD-x-EGACD-x(2)-

ELIVMFY3(4)-x(2)-G.

NAME: Cyclic nucleotide-binding domain signature 2.

CONSENSUS: ELIVMFI-G-E-x-EGASI-ELIVMI-x(5,11)-R-ESTAQI-A-x-

55 ELIVMAD-x-ESTACVD.

NAME: camp/cGMP binding motif.

NAME: EF-hand calcium-binding domain.

CONSENSUS: D-x-EDNSJ-{ILVFYW}-EDENSTGJ-EDNQGHRKJ-{GP}-

ELIVMCD-EDENQSTAGCD-x(2)-

CONSENSUS: EDED-ELIVMFYWD.

5

NAME: Actinin-type actin-binding domain signature L.

CONSENSUS: EEQD-x(2)-EYJ-x(2)-W-x-N.

NAME: Actinin-type actin-binding domain signature 2.

- LDASHQD-x-EDAZJ-ENVIJ-ENDZ-x-ENVIJ-ENZGJ-x-ENVIJ-

ELIVMD-x-EDEAGD-x(4)-

CONSENSUS: ELIVMI-x-ELMI-ESAGI-ELIVMI-ELIVMID-W-x-ELIVMI(2).

NAME: Anaphylatoxin domain signature.

15 CONSENSUS: ECSHI-C-x(2)-EGAPI-x(7,8)-EGASTDEQRI-C-EGASTDEQLI-

-(S)x-ENGAGTZABJ-(P,E)x

CONZENZUZ: ECE3-x(P'3)-C-C.

NAME: Anaphylatoxin domain profile.

20

NAME: Apple domain.

CONZENSUS: C-x(3)-ELIVMFYJ-x(5)-ELIVMFYJ-x(3)-EDENQJ-

 $\mathbb{L} \mathbb{L} \mathbb{I} \mathbb{V} \mathbb{M} \mathbb{F} \mathbb{V} \mathbb{J} - \times (\mathbb{J} \mathbb{D}) - \mathbb{C} - \times (\mathbb{J}) - \mathbb{C} - \mathbb{T} - \mathbb{C}$

CONSENSUS: x(4)-C-x-ELIVMFYJ-F-x-EYJJ-x(JJ-J4)-C-x-ELIVMFYJ-

25 ERKI-x-ESTI-x(14-15)-

CONSENSUS: S-G-x-E2TJ-ELIVMFYJ-x(2)-C.

NAME: Band 4.1 family domain signature 1.

CONSENSIONS: M-ELIVI-x(3)-EKRGI-x-ELIVII-x(2)-EGHI-x(0-2)-

30 ELIVMFI-x(b-8)-ELIVMFI-

CONSENSUS: $x(3-5)-F-\mathbb{C}Y\mathbb{I}-x(2)-\mathbb{C}D\mathbb{E}NS\mathbb{I}$.

NAME: Band 4.1 family domain signature 2.

- CONSENSUS: CHYUJ-x(P)-EVT29N3QJ-LVT29N3QD-x(P)-EUYHJ-x(P)-

35 $\mathbb{E}ACVJ-x(2)-\mathbb{E}LMJ-x(2)-$

NAME: Band 4.1 family domain profile.

40 NAME: Clq domain signature.

CONSENSUS: F-x(5)-ENDJ-x(4)-EFYULJ-x(6)-F-x(5)-G-x-Y-x-F-x-

[FY].

NAME: C-terminal cystine knot signature.

45 CONSENSUS: C-C-x(l3)-C-x(2)-EGNJ-x(l2)-C-x-C-x(2,4)-C.

NAME: C-terminal cystine knot profile.

NAME: CUB domain profile.

50 NAME:

Death domain profile.

NAME: EGF-like domain signature 1.

CONSENSUS: C-x-C-x(5)-G-x(2)-C.

55

NAME: EGF-like domain signature 2.

CONSENSUS: $C-x-C-x(2)-\mathbb{E}GP\mathbb{J}-\mathbb{E}FY\mathbb{J}-x(4-A)-C$.

NAME: Calcium-binding EGF-like domain pattern signature. CONSENSUS: EDEQNI-x-EDEQNI(2)-C-x(3,14)-C-x(3,7)-C-x-EDNI-x(4)-EFYI-x-C.

5 NAME: Laminin-type EGF-like (LE) domain signature. CONSENSUS: C-x(1,2)-C-x(5)-G-x(2)-C-x(2)-C-x(3,4)-EFYW1-x(3,15)-C.

NAME: Coagulation factors 5/8 type C domain (FAS8C)

15 NAME: Coagulation factors 5/8 type C domain (FAS8C) signature 2.

CONSENSUS: P-x(8-10)-ELMD-R-x-EGEB-ELIVPB-x-G-C.

NAME: Forkhead-associated (FHA) domain profile.

NAME: Fibrinogen beta and gamma chains C-terminal domain signature.

CONSENSUS: W-W-ELIVMFYWJ-x(2)-C-x(2)-EGSAJ-x(2)-N-G.

25 NAME: Type I fibronectin domainCONSENSUS: C-x(b-8)-ELFYI-x(5)-EFYWI-x-ERKI-x(8-10)-C-x-Cx(b-9)-C-

NAME: Type II fibronectin collagen-binding domainCONSENSUS: C-x(2)-P-F-x-EFYWID-x(7)-C-x(8,10)-W-C-x(4)EDNSR3-EFYWID-x(3,5)-EFYWID-xCONSENSUS: EFYWID-C.

NAME: Hemopexin domain signature.
CONSENSUS: ELIFATI-x(3)-W-x(2,3)-EPEI-x(2)-ELIVMFYI-EDENGSIESTAI-EAVI-ELIVMFYI.

NAME: Kringle domain signature.
CONSENSUS: EFYI-C-R-N-P-EDNRI.

NAME: Kringle domain profile.

NAME: LDL-receptor class A (LDLRA) domain signature.

CONSENSUS: C-EVILMAI-x(5)-C-EDNHI-x(3)-EDENGHTI-C-x(3,4)-

45 CSTADED-EDEHD-EDED-x(1,5)CONSENSUS: C.

35

40

NAME: LDL-receptor class A (LDLRA) domain profile.

50 NAME: C-type lectin domain signature.

CONSENSUS: C-ELIVMFYATGI-x(5,12)-EWLI-x-EDNSRI-x(2)-C-x(5,6)
EFYWLIVSTAI-ELIVMSTAI
CONSENSUS: C.

55 NAME: C-type lectin domain profile.

NAME: Link domain signature-CONSENSUS: C-x(15)-A-x(3)-C

NAME: Osteonectin domain signature 1-

CONSENSUS: C-x-EDND-x(2)-C-x(2)-G-EKRHD-x-C-x(6,7)-P-x-C-x-C-x(3,5)-C-P.

x(3,5)-

NAME: Osteonectin domain signature 2. CONSENSUS: F-P-x-R-EIMJ-x-D-W-L-x-ENQJ.

NAME: Somatomedin B domain signature.

10 CONSENSUS: C-x-C-x(3)-C-x(5)-C-C-x-EDNJ-EFYJ-x(3)-C.

NAME: Thyroglobulin type-1 repeat signature.

CONSENSUS: $\mathbb{C}FYUHPJ-x-P-x-C-x(3,4)-G-x-\mathbb{C}FYUJ-x(3)-Q-C-x(4,10)-$

C-EFYW3-C-V-x(3-4)~

15 CONSENSUS: ESG1.

NAME: P-type 'Trefoil' domain signature.

 $-2-2-(1)\times-2-(1)\times-1$

EFYWHI.

20

NAME: Cellulose-binding domain, bacterial type.

CONSENSUS: W-N-ESTAGRI-ESTDNI-ELIVMI-x(2)-EGSTI-x-EGSTI-x(2)-

ELIVMFT3-EGAD.

25 NAME: Cellulose-binding domain, fungal type.

C-G-HMJ-(B-E)x-D-(B)x-D-(E)x-D-(F-F)x-D-(F-F)

 $\mathbb{C}FYUMJI-x(2)-Q-C$

NAME: Chitin recognition or binding domain signature.

30 CONZENZUZ: C-x(4,5)-C-C-2-x(2)-G-x-C-G-x(4)-EFYUI-C.

NAME: Barwin domain signature 1.

CONSENSUS: C-G-EKRI-C-L-x-V-x-N.

35 NAME: Barwin domain signature 2.

CONSENSUS: V-EDND-Y-EEQJ-F-V-EDND-C.

NAME: BIR repeat.

CONZENZUZ: ENATZJ-(LVZJ-Z)-R-x(3-7)-EVMJ-x(LL-LVZJ-Z)-ECNATZ-G-

40 CLMFJ-X-CFYHDAJ-X(4)-

CONZENZUZ: $\mathbb{L}DEZL\mathbb{I}-X(5^{-3})-C-X(5)-C-X(7)-\mathbb{L}M\mathbb{I}-X(7)-H-X(4)-$

EPRSDJ-X-C-X(2)-ELIVMAJ.

NAME: WAP-type 'four-disulfide core' domain signature.

45 CONSENSUS: $C-x-\{C\}-EDNJ-x(2)-C-x(5)-C-C$.

NAME: Phorbol esters / diacylglycerol binding domain.

CONZENZUZ: H-x-ELIVMFYWJ-x(8-11)-C-x(2)-C-x(3)-ELIVMFCJ-

x(5,10)-C-x(2)-C-x(4)-EHD3-

50 CONZENZUZ: x(2)-C-x(5,7)-C.

NAME: C2 domain signature.

CONZENZUZ: EACGI-x(2)-L-x(2,3)-D-x(1,2)-ENGSTLIFI-EGTMRI-x-

ESTAPU-D-EPAU-EFYU.

NAME: C2-domain profile.

NAME: CAP-Gly domain signature.

CONSENSUS: G-x(B-1D)-EFYWD-x-G-ELIVMD-x-ELIVMFYD-x(4)-G-K-

-D-K-ERATZI-S-x-EHMI

CONSENSUS: x(2)-ELY3-F.

5 NAME: Ly-6 / u-PAR domain signature.

CONSENSUS: EEQRI-C-ELIVMFYAHI-x-C-x(5-a)-C-x(3-a)-EEDNQSTVI-

 $C-\{C\}-x(5)-C-$

CONSENSUS: x(75-54)-C.

10 NAME: MAM domain signature.

CONSENSUS: G-x-ELIVMFYJ(2)-x(3)-ECTJJ-x(JD-1JJ)-ELVJ-x(4)-

ELIVMFD-x(6-7)-C-ELIVMD-x-

CONSENSUS: $F-x-\mathbb{E}LIVMFYI-x(3)-\mathbb{E}GSCI$.

15 NAME: MAM domain profile.

NAME: PH domain profile.

NAME: Phosphotyrosine interaction domain (PID) profile.

20 NAME: Src homology 2 (SH2) domain profile.

NAME: Src homology 3 (SH3) domain profile.

25 NAME: VWFC domain signature.

->-(2123)-->-(4124)-(4124)-->->-(2131370)---

x(9,36)-C-C-x(2,4)-C.

NAME: WW/rsp5/WWP domain signature.

30 CONSENSUS: W-x(P,zll)-EYYJ-EYYJ-x(b,z)-EGTTRJ-EGSTRCRJ-

EFYU3-x(2)-P.

NAME: WW/rsp5/WWP domain profile.

35 NAME: ZP domain signature.

CONSENSUS: LLIVMFYWJ-x(7)-ESTANDUJ-x(3)-ELIVMFYWJ-x-

ELIVMFYUD-x-ELIVMFYUD-x(2)-C-

CONSENSUS: ELIVMFYWD-x-ESTD-EPSLD-x(2,4)-EDENSD-x-ESTADNQLFD-

x(L)-ELIVM3(2)-x(3-4)-

40 CONSENSUS: C-

NAME: S-layer homology domain signature.

CONSENSUS: CLVFYT3-x-CDA3-x(2,5)-CDNGSATPHY3-CWYFPDA3-x(4)-

ELIV3-x(2)-EGTALV3-

45 CONSENSUS: x(4,6)-CLIVFYCJ-x(2)-G-x-EQGTAJ-x(2,3)-EMFYAJ-x-

EPGAVI-x(3,10)-CLIVMAI-

CONSENSUS: ESTKRI-ERYI-x-EEQI-x-ESTALIVMI.

NAME: 'Homeobox' domain signature.

50 CONSENSUS: ELIVMTYAL-EASLVRI-x(2)-ELIVMSTACNI-x-ELIVMI-x(4)-

ELIVI-ERKNQESTAIYI-

CONSENSUS: ELIVFSTNKHI-W-EFYVCI-x-ENDQTAHI-x(5)-ERKNAIMUI.

NAME: 'Homeobox' domain profile.

55 NAME: 'Homeobox' antennapedia-type protein signature.

CONSENSUS: CLIVMFED-EFYD-P-W-M-EKRQTAD.

NAME: 'Homeobox' engrailed-type protein signature.

CONSENSUS: L-M-A-Q-G-L-Y-N.

NAME: 'Paired box' domain signature.

5 CONSENSUS: R-P-C-x(ll)-C-V-S.

NAME: 'POU' domain signature 1.

CONSENSUS: ERKQD-R-ELIMD-x-ELFD-G-ELIVMFYD-x-Q-x-EDNQD-V-G.

10 NAME: 'POU' domain signature 2.

CONSENSUS: S-Q-ESTJ-ETAJ-I-ESCJ-R-F-E-x-ELSQJ-x-ELIJ-ESTJ.

NAME: Zinc finger, C2H2 type, domain.

-H-(2,E)x-H-(B)x-E)UY7MVIJI-(E)x-D-(4,5)x-D

15

35

NAME: Zinc finger, C3HC4 type (RING finger), signature.

CONSENSUS: C-x-H-x-ELIVMFYD-C-x(2)-C-ELIVMYAD.

NAME: Nuclear hormones receptors DNA-binding region

20 signature.

CONSENSUS: C-x(2)-C-x-EDE3-x(5)-EHN3-EFY3-x(4)-C-x(2)-C-x(2)-

F-F-x-R.

NAME: GATA-type zinc finger domain-

25 CONSENSUS: C-x-EDNJ-C-x(4,5)-ESTJ-x(2)-W-EHRJ-ERKJ-x(3)-EGNJ-

-)-[ZA]-N-D-(4,E)x

NAME: Poly(ADP-ribose) polymerase zinc finger domain

signature.

30 CONSENSUS: C-EKRI-x-C-x(3)-I-x-K-x(3)-ERGI-x(16-18)-W-EFYHI-

H-x(2)-C.

NAME: Poly(ADP-ribose) polymerase zinc finger domain

profile.

NAME: Fungal Zn(2)-Cys(b) binuclear cluster domain

signature.

CONSENSUS: EGASTPVI-C-x(2)-C-ERKHSTACWI-x(2)-ERKHQI-x(2)-C-

x(5,12)-C-x(2)-C-x(6,1)-

40 · CONSENSUS: C.

NAME: Fungal Zn(2)-Cys(b) binuclear cluster domain profile.

NAME: Prokaryotic dksA/traR C4-type zinc finger.

45 CONSENSUS: $C-\mathbb{E} Z = -x - C - x(3) - I - x(3) - R - x(4) - C - x(2) - C$.

NAME: Copper-fist domain signature.

CONSENSUS: M-ELIVMFJ(3)-x(3)-K-EMYJ-A-C-x(2)-C-I-EKRJ-x-H-

EKR3-x(3)-C-x-H-x(8)-

50 CONSENSUS: EKRI-x-EKRI-G-R-P.

NAME: Copper fist DNA binding domain profile.

NAME: Leucine zipper pattern.

55 CONSENSUS: L-x(b)-L-x(b)-L-x(b)-L.

NAME: bZIP transcription factors basic domain signature.

CONSENSUS: EKNJ-x(1,3)-ERKZNAJ-x-(2)-ESAQJ(2)-x-ERXJ-x-(LONGNATNAJ-x-ERKJ).

NAME: Myb DNA-binding domain repeat signature 1.

5 CONSENSUS: W-EST3-x(2)-E-EDE3-x(2)-ELIV3.

NAME: Myb DNA-binding domain repeat signature 2.

CONSENSUS: W-x(2)-ELIJ-ESAGJ-x(4-5)-R-x(8)-EYWJ-x(3)-ELIVMJ-

10 NAME: Myc-type, 'helix-loop-helix' dimerization domain

signature.

CONSENSUS: CDENSTAPI-K-CLIVIWAGSNI-4FYWCPHKR}-CLIVTI-CLIVI-

-x-EDAT2MVIJJ-EVAT2J-(2)x

CONSENSUS: EVMFYHJ-ELIVMTAJ-{P}-{P}-ELIVMSRJ.

NAME: p53 tumor antigen signature.
CONSENSUS: M-C-N-S-S-C-M-G-G-M-N-R-R.

NAME: CBF-A/NF-YB subunit signature.

20 CONSENSUS: C-V-S-E-x-I-S-F-ELĨVMJ-T-ESGJ-E-A-ESCJ-EDEJ-EKRQJ-

NAME: CBF-B/NF-YA subunit signature.

CONSENSUS: $Y-V-N-A-K-Q-Y-x-R-\tilde{I}-L-K-R-R-x-A-R-A-K-L-E$.

25
NAME: 'Cold-shock' DNA-binding domain signature.

CONSENSUS: EFYJ-G-F-I-x(b,7)-EDERJ-ELIVMJ-F-x-H-x-ESTKRJ-x-

ELIVMFYD.

30 NAME: CTF/NF-I signature.

CONSENSUS: R-K-R-K-Y-F-K-K-H-E-K-R.

NAME: Ets-domain signature 1.

CONSENSUS: L-EFYWD-E@EDHD-F-ELID-ELV@KD-x-ELID-L.

35 NAME: Ets-domain signature 2.

CONSENSUS: CRKHJ-x(2)-M-x-Y-EDENQJ-x-ELIVMJ-ESTAGJ-R-ESTAGJ-

ELII-R-x-Y.

40 NAME: Ets-domain profile.

NAME: Fork head domain signature 1.

CONSENSUS: EKRD-P-EPTQD-EFYLVQHD-S-EFYD-x(2)-ELIVMD-x(3,4)-

EACD-ELIMD.

50

55

45

NAME: Fork head domain signature 2.

CONSENSUS: W-EQKRI-ENSI-S-ELIVI-R-H.

NAME: Fork head domain profile.

NAME: HSF-type DNA-binding domain signature.

CONSENSUS: L-x(3)-EYJ-K-H-x-N-x-EZJANJ-S-F-ELIVMJ-R-Q-L-

ENHU-x-Y-x-EFYWU-ERKHU-K-

CONZENSUS: ELIVMI.

NAME: Tryptophan pentad repeat (IRF family) signature.

CONSENSUS: W-x-EDNHI-x(5)-ELIVFI-x-EIVI-P-W-x-H-x(9,10)-EDEI-

x(2)-ELIVFI-F-EKRQI-x-

CONZENZUZ:

EURI-A.

NAME:

LIM domain signature.

NAME: LIN COMMIN Signature.

CONZENZUZ: C-x(5)-C-x(J2-5J)-EFYMHJ-H-x(5)-ECHJ-x(5)-C-x(5)-

5 C-x(3)-ELIVMF3.

NAME: LIM domain profile.

NAME: NF-kappa-B/Rel/dorsal domain signature.

10 CONSENSUS: F-R-Y-x-C-E-G.

NAME: MADS-box domain signature.

CONSENSUS: R-x-ERKJ-x(5)-I-x-EDNJ-x(3)-EKRJ-x(2)-T-EFYJ-x-

ERKI(3)-x(2)-ELIVMI-x-

15 CONSENSUS: K(2)-A-x-E-ELIVMI-ESTI-x-L-x(4)-ELIVMI-x-

ELIVM3(3)-x(b)-ELIVMF3-x(2)-

CONSENSUS: EFY3.

NAME: MADS-box domain profile-

20

NAME: T-box domain signature 1.

CONSENSUS: L-W-x(2)-EFCJ-x(3-4)-ENTJ-E-M-ELIVJ(2)-T-x(2)-G-

ERGI-EKRQI.

25 NAME: T-box domain signature 2.

CONZENZUZ: ELIVMYUJ-H-EPADHJ-EDENJ-EGSJ-x(3)-G-x(2)-W-M-x(3)-

EIVAI-x-F.

NAME: TEA domain signature.

30 CONSENSUS: G-R-N-E-L-I-x(2)-Y-I-x(3)-ETCI-x(3)-R-T-ERKI(2)-Q-

ELIVIJ-2-2-H-ELIVIJ-Q-V.

NAME: Transcription factor TFIIB repeat signature.

35 CONSENSUS: G-EKRU-x(3)-ESTADNU-x-ELIVMYAU-EĞSTAD(2)-ECVAD-

CLIVMD-CLIVMFYD-CLIVMAD-

CONSENSUS: EGSAD-ESTACD-

NAME: Transcription factor TFIID repeat signature.

40 CONSENSUS: Y-x-P-x(2)-EIFI-x(2)-ELIVII(2)-x-EXRHII-x(3)-P-

ERKQ3-x(3)-L-ELIVM3-F-x-

CONSENSUS: ESTND-G-EKRD-ELIVMD-x(3)-G-ETAGLD-EKRD-x(7)-EAGCD-

x(7)-ELIVMB.

45 NAME: TFIIS zinc ribbon domain signature.

CNSENSUS: C-x(2)-C-x(7)-ILJACABALA-EQHJ-ESACRJ-

x-EDEI-EDETI-EPGSEAI-

CONSENSUS: $x(b)-C-x(2-5)-C-x(3)-\mathbb{E}FU$.

50 NAME: TSC-22 / dip / bun family signature.

CONSENSUS: M-D-L-V-K-x-H-L-x(2)-A-V-R-E-E-V-E.

NAME: Prokaryotic transcription elongation factors signature

1.

x(2)-EIVI-x(3)-ELIVI-

CONSENSUS: $x(b)-G-D-x(2)-E-N-\mathbb{E}GSAJ-x-Y$.

-399-

NAME: Prokaryotic transcription elongation factors signature

2.

CONSENSUS: S-x(2)-S-P-ELIVMI-EAGI-x-ESAGI-ELIVMI-ELIVMYI-

x(4)-EDGI-EDEI.

5

NAME: DEAD-box subfamily ATP-dependent helicases signature.
CONSENSUS: ELIVMF1(2)-D-E-A-D-ERKEND-x-ELIVMFYGSTND.

NAME: DEAH-box subfamily ATP-dependent helicases signature.

10 CONSENSUS: EGSAHI-x-ELIVMFI(3)-D-E-EALIVI-H-ENECRI-

NAME: Eukaryotic putative RNA-binding region RNP-1

signature.

CONSENSUS: ERKI-G-{EDRKHPCG}-EAGSCIJ-EFYJ-ELIVAJ-x-EFYLMJ.

15

NAME: Fibrillarin signature.

CONSENSUS: EGZTJ-ELIVĀAPJ-V-Y-A-EVJJ-E-EFYJ-EZAJ-x-R-x(2)-R-

EDEI.

20 NAME: MCM family signature.

CONSENSUS: G-EIVIJ-ELVACJ(2)-EIVIJ-D-EDEJ-EFLJ-EDNSTJ.

NAME: MCM family domain.

25 NAME: XPA protein signature 1.

CONSENSUS: $C-x-\mathbb{E}DEJ-C-x(3)-\mathbb{E}LIVMFJ-x(J-2)-D-x(2)-L-x(3)-F-$

x(4)-C-x(2)-C.

NAME: XPA protein signature 2.

30 CONSENSUS: ELIVMJ(2)-T-EKRJ-T-E-x-K-x-EDEJ-Y-ELIVMFJ(2)-x-D-

x-EDEI-

NAME: XPG protein signature 1-

CONSENSUS: EVID-EKRED-P-x-EFYILD-V-F-D-G-x(2)-EPILD-x-ELVCD-

35 K•

40

NAME: XPG protein signature 2.

CONSENSUS: EGSD-ELIVMD-EPERD-EFYSD-ELIVMD-x-A-P-x-E-A-EDED-

EPASI-EQSI-ECLMI.

NAME: Bacterial regulatory proteins, anaC family signature.

CONSENSUS: EKRQJ-ELIVMAJ-x(2)-EGSTALIVJ-{FYWPGDN}-x(2)-

ELIVMSAD-x(4,9)-ELIVMFD-

CONSENSUS: x(2)-ELIVMSTAD-EATADIDATCED-x(3)-EATADIDECTIONSTAD-ELIVMYYD-

45 x(4-5)-EEFY3-x(3)-

CONSENSUS: EFYIVAI-(F)-EAVIND-x-ENSHAPKLI-x-ENSHAPKLI-

EPARLI.

NAME: Bacterial regulatory proteins, araC family DNA-binding

50 domain profile.

NAME: Bacterial regulatory proteins, arsR family signature.

CONSENSUS: C-x(2)-D-ELIVMI-x(6)-ESTI-x(4)-S-EHYRI-EHQI.

55 NAME: Bacterial regulatory proteins, asnC family signature.

-EYAMVIJI-(5)x-EAZDI-EMVIJI-EANGI-(5)x-EQATZDI : ZUZNIZNO)

EGND-ELIVMSTD-ESTD-x(b)-R-

. CONSENSUS: ELVTD-x(2)-ELIVMD-x(3)-G.

NAME: Bacterial regulatory proteins, crp family signature.

CONSENSUS: CLIVMI-CSAT-CENVHID-CGAT-x-CLIVMYAI
ELIVACI-CGAT-x-CSAT

5 CONZENSUS: x(2)-EMTJ-x-EGTTJ-x-ELIVMFJ-x(2)-ELIVMFJ.

NAME: Bacterial regulatory proteins, deoR family signature.

CONSENSUS: R-x(3)-ELIVM3-x(3)-ELIVM3-x(16,17)-ESTA3-x(2)-TELIVM3-ERH3-EKMN3-D-

10 CONSENSUS: CLIVMFI.

NAME: Bacterial regulatory proteins, gntR family signature.

CONSENSUS: CLIVAPKRI-EPILVI-x-CEQTIVMRI-x(2)-CLIVMI-x(3)
CLIVMFYKI-x-CLIVFII-

15 CONSENSUS: ELDMGSTKJ-ERGTLVJ-x-ESTAIVPJ-ELIVAJ-x(2)-ESTAGVJ-ELIVMFYHJ-x(2)-ELMAJ.

20 x(2)-ELIVM3-EFYH3-EDN3.

25 CONSENSUS: ELIVMFYAND-ELIVMCD.

NAME: Bacterial regulatory proteins, luxR family signature.

CONSENSUS: EGDCI-x(2)-ENSTAVYI-x(2)-EIVI-EGSTAI-x(2)
ELIVMFYWCTI-x-ELIVMFYWCRI-x(3)-

30 CONSENSUS: ENST3-ELIVM3-x(5)-ENRHSA3-ELIVMSTA3-x(2)-EKR3.

NAME: Bacterial regulatory proteins, lysR family signature.
CONSENSUS: ENGKRHSTAGI-ELIVMFYTAI-x(2)-ESTAGLVI-ESTAGI-x(4)ELIVMYCTGRI-EPSTANLVERI-

35 CONSENSUS: x-EPSTAGQVI-EPSTAGAVNFI-ELIVMFAI-ESTAGHI-x(2)-ELIVMFI-x(2)-ELIVMFI-x(2)-ELIVMFI-x(2)-ELIVMFI-x(3)-ELIVMFI-x(

NAME: Bacterial regulatory proteins, mark family signature.
40 CONSENSUS: ESTNAI-ELIAI-x-ERNGSI-x(4)-ELMI-EEIVI-x(2)-EGESI-ELFYWI-ELIVCI-x(7)CONSENSUS: EDNI-ERKQGI-ERKI-x(b)-T-x(2)-EGAI.

NAME: Bacterial regulatory proteins, merR family signature.
45 CONSENSUS: EGSAD-x-ELIVMFAD-EASMD-x(2)-ESTACLIVD-EGSDENQRDELIVCD-ESTANHKD-x(3)CONSENSUS: ELIVMD-ERHFD-x-EYWD-EDEQD-x(2,3)-EGHDNQDELIVMFD(2).

50 NAME: Bacterial regulatory proteins, tetR family signature.

CONSENSUS: G-ELIVMFYSD-x(2,3)-ETSD-ELIVMTD-x(2)-ELIVMD-x(5)
ELIVQSD-ESTAGENQHD-x
CONSENSUS: EGPARD-x-ELIVMFD-EFYSTD-x-EHFYD-EFVD-x-EDNSTD-K
x(2)-ELIVMD.

NAME: Transcriptional antiterminators bglG family signature.

CONSENSUS: ESTI-x-H-x(2)-EFAI(2)-ELIVNI-EEQKI-R-x(2)-EQNKI.

Sigma-54 factors family signature 1.

CONSENSUS: P-ELIVM3-x-ELIVM3-x(2)-ELIVM3-A-x(2)-ELIVMF3-x(2)-

EHSD-x-S-T-ELIVMD-S-R.

Sigma-54 factors family signature 2.

CONZENZUZ: R-R-T-EIVI-EATI-K-Y-R.

Sigma-54 factors family profile. NAME:

10 NAME: Sigma-70 factors family signature 1.

CONSENSUS: EDED-ELIVMFD(2)-EHEQSD-x-G-x-ELIVMFAD-G-L-

CLIVMFYED-x-EGSAMD-CLIVMAPD.

NAME: Sigma-70 factors family signature 2.

ESTND-x(2)-EDEQD-ELIVMD-EGASD-x(4)-ELIVMFD-EPSTGD-15 CONZENZUZ:

x(3)-ELIVMAI-x-ENGRI-

CONZENSUS: ELIVMAD-EEQHD-x(3)-ELIVMFUD-x(2)-ELIVMD.

Sigma-70 factors ECF subfamily signature. NAME:

20 CONZENZUZ: USTAIVD-CPQDELD-CDED-CLIVD-CLIVAD-Q-x-CSTAVD-

ELIVMFYCD-ELIVMAKD-x-

CONZENZUZ: EGSTAIV3-ELIMFYUQ3-x(12,14)-ESTAP3-EFYU3-ELIF3-

x(2)-[[V]].

25 NAME: Sigma-54 interaction domain ATP-binding region A

signature.

CONZENSUS: ELIVMFYD(3)-x-G-EDEQD-ESTED-G-ESTAVD-G-K-x(2)-

CLIVMFYJ.

30 NAME: Sigma-54 interaction domain ATP-binding region B

signature.

CONZENZUZ: -CASSI-x-ELIVMFI-x(2)-A-EDNEQASHI-EGNEKI-G-ESTIMI-

CLIVMFY3(3)-EDE3-CEK3-

CONZENZUZ: ELIVMI.

35

NAME: Sigma-54 interaction domain C-terminal part signature.

CONSENSUS: EFYWD-P-EGSD-N-ELIVMD-R-EEQD-L-x-ENHATD.

NAME: Sigma-54 interaction domain profile.

40

NAME: Single-strand binding protein family signature 1.

CLIVMF3-CNST3-CKRT3-CLIVM3-x-CLIVMF3(2)-G-ENHRK3-CONZENZUZ:

ELIVMD-EGSTD-x-EDETD.

45 Single-strand binding protein family signature 2. NAMF:

T-x-W-EHYJ-ERNSJ-ELIVMJ-x-ELIVMFJ-EFYJ-ENGKRJ. **CONSENSUS:**

Bacterial histone-like DNA-binding proteins signature. NAMF:

CONZENZUZ: EGSKI-F-x(2)-ELIVMFI-x(4)-ERKEQAI-x(2)-ERSTI-x-

50 EGAI-x-EKNI-P-x-T.

> NAME: Dps protein family signature 1.

CONZENZUZ: $H-\mathbb{E}FUJ-x-\mathbb{E}LIVMJ-x-G-x(5)-\mathbb{E}LVJ-H-x(3)-\mathbb{E}DEJ$

55 NAME: Dps protein family signature 2.

CONZENSUS: LLIVMFYD-EDHD-x-ELIVMD-EGAD-E-R-x(3)-ELIFD-EGDND-

 \times (2)-EPAI.

NAME: DNA repair protein radC family signature.

CONSENSUS: H-N-H-P-S-G.

NAME: recA signature.

5 CONSENSUS: A-L-EKRD-EIFD-EFYD-ESTADD-ELIVMQD-R.

NAME: RecF protein signature 1.

EFYD-ELIVMD-D.

10

NAME: RecF protein signature 2.

CONSENSUS: ELIVMFYD(2)-x-D-x(2,3)-EEAD-EEHD-L-D-x(2)-EKRHD-

·x(3)-L.

15 NAME: RecR protein signature.

CONSENSUS: C-x(2)-C-x(3)-C-x-1-C-x(4)-C-x-1-C-x(4)-R.

NAME: Histone H2A signature.
CONSENSUS: EACI-G-L-x-F-P-V.

20

NAME: Histone H2B signature.

CONSENSUS: EKRJ-E-ELIVMJ-EEQJ-T-x(2)-EKRJ-x-ELIVMJ(2)-x-...

EPAGD-EDED-L-x-EKRD-H-A-

CONSENSUS: ELIVMI-ESTAI-E-G.

25

NAME: Histone H3 signature 1.

CONSENSUS: K-A-P-R-K-Q-L.

NAME: Histone H3 signature 2.

30 CONSENSUS: P-F-x-ERAI-L-EVAI-EKRQI-EDEGI-EIVI-

NAME: Histone H4 signature.

CONSENSUS: G-A-K-R-H.

35 NAME: HMG1/2-signature.

NAME: HMG-I and HMG-Y DNA-binding domain (A+T-hook).
CONSENSUS: EAT3-x(1,2)-ERK3(2)-EGP3-R-G-R-P-ERK3-x.

40

NAME: HMG14 and HMG17 signature.
CONSENSUS: R-R-S-A-R-L-S-A-ERKI-P.

NAME: Bromodomain signature.

45 CONSENSUS: ESTANVED-x(2)-F-x(4)-EDNSD-x(5,7)-EDENQTFD-Y-

EHFYJ-x(2)-ELIVMFYJ-x(3)-

CDNS=NZUS: CENVIJU - x (4) - CENVIJU - x (6, 8) - Y - x (12, 13) - (2) -

N-ESACFJ-x(2)-EFYJ.

50 NAME: Bromodomain profile.

NAME: Chromo domain signature-

CONSENSUS: EFYLJ-x-ELIVMCJ-EKRJ-W-x-EGDNRJ-EFYWLEJ-x(5,6)-

-(E)x-UNGTZqJ-U-UZJJ-U-UTZJ

55 CONSENSUS: ELIVMCI.

NAME: Chromo and chromo shadow domain profile.

NAME: Regulator of chromosome condensation (RCCL) signature

լ.

CONSENSUS: G-x-N-D-x(2)-EAVJ-L-G-R-x-T.

5 NAME: Regulator of chromosome condensation (RCCL) signature

2.

CONSENSUS: ELIVMFAD-ESTAGOD(2)-G-x(2)-H-ESTAGLID-ELIVMFAD-x-

ELIVMD.

10 NAME: Protamine Pl signature.

-S-x-Z-x-LTZJ-(2,3)x-R-LYJNJ-R-LVJU : 2UZNJZNO)

NAME: Nuclear transition protein 1 signature.

CONSENSUS: S-K-R-K-Y-R-K.

15

45

NAME: Nuclear transition protein 2 signature 1.

CONSENSUS: H-x(3)-H-S-ENSJ-S-x-P-Q-S.

NAME: Nuclear transition protein 2 signature 2.

20 CONSENSUS: K-x-R-K-x(2)-E-G-K-x(2)-K-EKRI-K.

NAME: Ribosomal protein Ll signature-

CONSENSUS: EIMI-x(2)-ELVAI-x(2-3)-ELVAI-x(2)

EGSNHJ-EPTKRJ-EKRAVJ-G-x-

25 CONSENSUS: ELMF3-P-EDENSTK3.

NAME: Ribosomal protein L2 signature.

CONSENSUS: P-x(2)-R-G-ESTAIVJ(2)-x-N-EAPKJ-x-EDEJ.

30 NAME: Ribosomal protein L3 signature.

CONSENSUS: EFLJ-x(b)-EDNJ-x(2)-EAGSJ-x-ESTJ-x-G-EKRHJ-G-x(2)-

G = x(3) - R

NAME: Ribosomal protein L5 signature-

35 CONSENSUS: ELIVMI-x(2)-ELIVMI-ESTACI-EGEI-EQVI-x(2)-ELIVMAI-

x-ESTCD-x-ESTAGD-EKRD-CONSENSUS: x-ESTAD.

NAME: Ribosomal protein Lb signature 1.

40 CONSENSUS: EPSI-EDENSI-x-Y-K-EGAI-K-G-ELIVMI.

NAME: Ribosomal protein Lb signature 2.

CONSENSUS: Q-x(3)-ELIVMJ-y-(2)-EKRJ-x(2)-R-x-F-x-D-G-ELIVMJ-Y-

ELIVMI-x(2)-EKRI.

NAME: Ribosomal protein L9 signature.

CONSENSUS: G-x(2)-EGN3-x(4)-V-x(2)-G-EFY3-x(2)-N-EFY3-L-x(5)-

-ENTZJ-(E)x-EADJ

50 NAME: Ribosomal protein LLD signature.

CONSENSUS: EDEHD-x(2)-EGSD-ELIVMFD-ESTND-EVAD-x-EDEQKD-

ELIVMAD-x(2)-ELIMD-R.

NAME: Ribosomal protein Lll signature.

55 CONSENSUS: ERKND-x-ELIVMD-x-G-ESTD-x(2)-ESNQD-ELIVMD-G-x(2)-

ELIVMI-x(0,1)-EDENGI-

NAME: Ribosomal protein Ll3 signature.

-404-

ENQEKRAJ-x(5)-ELIVMJ-x-EAIVJCONSENSUS: ELFYJ-x-EGDNJ.

5 NAME: Ribosomal protein Ll4 signature.
CONSENSUS: EGAl-ELIVI(3)-x(9,l0)-EDNSI-G-x(4)-EFYI-x(2)-ENTI-

x(2)-V-ELIVI.

NAME: Ribosomal protein L15 signature.

10 CONSENSUS: K-ELIVMJ(2)-EGALJ-x-EGTJ-x-ELIVMJ-x(2,5)-ELIVMJ-

x-ELIVMF3-x(3,4)-

. O-(E)x-EMVIJJ-(E)x-A-(S)x-ETZJ-EJ7MVIJJ :ZUZNJZNOJ

NAME: Ribosomal protein Llb signature 1.

15 CONSENSUS: EKRI-R-x-EGSACI-EKQVAI-ELIVMI-W-ELIVMI-EKRI-ELIVMI-ELFYI-EAPII-

NAME: Ribosomal protein Llb signature 2. CONSENSUS: R-M-G-x-EGRI-K-G-x(4)-EFWKRI.

20

NAME: Ribosomal protein Ll7 signature.

CONSENSUS: I-x-ESTI-EGTI-x(2)-EKRI-x-K-x(6)-EDEI-x-ELIMVIELIVMTI-T-x-ESTAGI-EKRI.

25 NAME: Ribosomal protein Ll9 signature.

CONSENSUS: ERTJ-EKRSVYJ-EGSAJ-x-V-ERSJ-EKRJ-ESAJ-K-L-Y-Y-L-R.

NAME: Ribosomal protein L2O signature.

CONSENSUS: K-x(3)-EKRCJ-x-ELIVMJ-W-EIVJ-ESTNALVJ-R-ELIVMJ-N-

30 x(3)-ERKHI-

35

NAME: Ribosomal protein L22 signature.

CONSENSUS: ERKQNJ-x(4)-ERHJ-EGASJ-x-G-EKRQSJ-x(9)-EHDNJELIVMJ-x-ELIVMSJ-x-ELIVMJ.

40 NAME: Ribosomal protein L23 signature.

CONSENSUS: ERKJ(2)-EAMJ-EIVFYTJ-EIVJ-ERKTJ-L-ESTANQKJ-x(7)
ELIVMFTJ.

NAME: Ribosomal protein L24 signature.

45 CONSENSUS: CGDENI-D-x-V-x-EIVI-ĒLIVMAI-x-G-x(2)-EKAI-EGNI-x(2,3)-EGAI-x-EIVI.

NAME: Ribosomal protein L27 signature.
CONSENSUS: G-x-ELIVMI(2)-x-R-Q-R-G-x(5)-G.

50

NAME: Ribosomal protein L27 signature.

CONSENSUS: EKNQSI-EPSTLI-x(2)-ELIMFAB-EKRGSANI-x-ELIVYSTAI
EKRI-EKRI-EKRHI-EDESTANRLI
CONSENSUS: ELIVI-A-EKRCQVTI-ELIVMAI.

55

CONSENSUS: x(10)-ELMS1-ELIV1-x(2)-ELIVA1-x(2)-ELMFY1-EIVT1.

NAME: Ribosomal protein L31 signature.

CONSENSUS: H-P-F-EFYJ-ETIJ-x(9)-G-R-EAVJ-x-EKRJ.

5 NAME: Ribosomal protein L33 signature

NAME: Ribosomal protein L33 signature.

CONSENSUS: Y-x-ETZJ-x-EKRJ-ENSJ-x(4)-EPATJ-x(1,2)-ELIVMJ-

TEAD-x(2)-K-TFYD-TCSDD.

10 NAME: Ribosomal protein L34 signature.

CONSENSUS: K-ERGD-T-EFYWLD-EEQSD-x(5)-EKRHSD-x(4,5)-G-F-x(2)-

R.

NAME: Ribosomal protein L35 signature.

15 CONSENSUS: ELIVMD-K-ETVD-x(2)-EGSAD-ESAILD-x-K-R-ELIVMFYD-

EKRLI.

NAME: Ribosomal protein L36 signature.

-x-EMVIJa-x-ENMVIJa-x-R-x(3)-ENVIJa-x-EMVIJa-x-C-x(2)-C-x(2)-C-x(2)-C-x-C-x(3)-ENZUZNO3

20 C-x(3-4)-EKRJ-H-x-Q-x-Q-

NAME: Ribosomal protein Lle signature.

ESGAD-x(7)-ERKD-G-H.

25

NAME: Ribosomal protein Lbe signature.

CONSENSUS: N-x(2)-P-L-R-R-x(4)-EYJJ-V-I-A-T-S-x-K.

NAME: Ribosomal protein L7Ae signature.

30 CONSENSUS: ECAD-x(4)-EIVJ-P-EFYJ-x(2)-ELIVMJ-x-EGSQJ-EKRQJ-

x(2)-L-G.

NAME: Ribosomal protein LlDe signature.

CONSENSUS: R-x-A-EfyWJ-G-K-EPJJ-x-G-x(2)-A-R-V.

35 -

NAME: Ribosomal protein Ll3e signaturé.

CONSENSUS: EKRI-Y-x(2)-K-ELIVMI-R-ESTAI-G-EKRI-G-F-ESTI-L-x-

Ε.

40 NAME: Ribosomal protein Ll5e signature-

CONSENSUS: EDEI-EKRI-A-R-x-L-G-EFYI-x-ESAPI-x(2)-G-

ELIVMFY3(4)-R-x-R-V-x-R-G.

NAME: Ribosomal protein Llae signature.

45 CONSENSUS: CKREJ-x-L-x(2)-CPSJ-CKRJ-x(2)-CRHJ-CPSAJ-x-CLIVMJ-

ENSD-ELIVMD-x-ERKD-

CONSENSUS: ELIVMI-

NAME: Ribosomal protein Ll9e signature.

50 CONSENSUS: R-x-EKRII-x(5)-EKRII-x(3)-EKRHI-x(2)-G-x-G-x-R-x-G-

x(3)-A-R-x(3)-EKQ3-

CONSENSUS: $\times (2) - U - \times (7) - R - \times (2) - L - \times (3) - R$

NAME: Ribosomal protein L2le signature.

55 CONSENSUS: G-EDED-x-V-x(10)-EGVD-x(2)-EFYHD-x(2)-EFYD-x-G-x-

T-G.

NAME: Ribosomal protein L24e signature.

NAME: Ribosomal protein L27e signature.
5 CONSENSUS: G-K-N-x-W-F-F-x-K-L-R-F>.

NAME: Ribosomal protein L3De signature l.

CONSENSUS: ESTAD-x(5)-G-x-EQKRD-x(2)-ELIVMD-EKQTD-x(2)-EKRD-

x-G-x(2)-K-x-ELIVM3(3)-

NAME: Ribosomal protein L3De signature 2.

CONSENSUS: EDEI-L-G-ESTAI-x(2)-G-EKRI-x(b)-ELIVMI-x-ELIVMI-x-

EDENI-x-G.

15 NAME: Ribosomal protein L3le signature.
CONSENSUS: V-EKRJ-ELIVMJ-x(3)-ELIVMJ-N-x-EAKJ-x-W-x-EKRJ-G.

NAME: Ribosomal protein L32e signature.

CONSENSUS: F-x-R-x(4)-EKRJ-x(2)-EKRJ-ELIVMJ-x(3)-U-R-EKRJ-

20 x(2)-6.

- NAME: Ribosomal protein L34e signatureCONSENSUS: Y-x-EST3-x-S-ENY3-x(5)-EKR3-T-P-G.

25 NAME: Ribosomal protein L35Ae signature.
CONSENSUS: G-K-ELIVM3-x-R-x-H-G-x(2)-G-x-V-x-A-x-F-x(3)-ELI3-P.

NAME: Ribosomal protein L3Le signature.

30 CONSENSUS: P-Y-E-EKRI-R-x-ELIVMI-EDEI-ELIVMI(2)-EKRI-

NAME: Ribosomal protein L37e signature-

CONSENSUS: $G-T-x-\mathbb{E}SAJ-x-G-x-\mathbb{E}KJ-x(3)-\mathbb{E}STJ-x(0-1)-H-x(2)-C-x-$

R-C-G-

NAME: Ribosomal protein L39e signature.

CONSENSUS: EKRAD-T-x(3)-ELIVMD-EKRQFD-x-ENHSD-x(3)-R-ENHYD-W-

R-R.

40 NAME: Ribosomal protein L44e signature.
CONSENSUS: K-x-ETVJ-K-K-x(2)-L-EKRJ-x(2)-C.

•

NAME: Ribosomal protein S2 signature l.
CONSENSUS: ELIVMFAD-x(2)-ELIVMFYCD(2)-x-ESTACD-EGSTANGEKRD-

45 ESTALVI-EHYI-ELIVMFI-G.

NAME: Ribosomal protein S2 signature 2.

CONSENSUS: P-x(2)-ELIVMFJ(2)-ELIVMSJ-x-EGDNJ-x(3)-EDENLJ-

x(3)-ELIVMB-x-E-x(4)-

50 CONSENSUS: EGNQKRH3-ELIVM3-EAP3.

NAME: Ribosomal protein S3 signature.

CONSENSUS: EGTAJ-EKRJ-x(L)-G-x-ELIVMTJ-x(2)-EHZSQAJ-EKRJ-x(1,3)-

ELIVFCAD-x(3)-ELIVD-

55 CONSENSUS: $\mathbb{C}DENQJ-x(7)-\mathbb{C}LMTJ-x(2)-G-x(2)-G$.

NAME: Ribosomal protein S4 signature.

CONSENSUS: ELIVMI-EDEI-x-R-L-x(3)-ELIVMCI-EVMFYHQI-EKRTI-

x(3)-ESTAGCFJ-x-ESTJ-x(3)-

CONSENSUS: ESAID-EKRD-x-ELIVMFD(2).

5 NAME: Ribosomal protein S5 signature.

CONSENSUS: G-EKRQJ-x(3)-EFYJ-x-EACVJ-x(2)-ELIVMAJ-ELIVMJ-

EAGU-EDNU-x(2)-G-x-

CONSENSUS: ELIVMI-G-x-ESAGI-x(5,6)-EDEQI-ELIVMI-x(2)-A-

ELIVMF3-

10

NAME: Ribosomal protein Sb signature.

CONSENSUS: G-x-EKRCJ-EDENQRHJ-L-ESAJ-Y-x-I-EKRNSAJ.

NAME: Ribosomal protein S7 signature.

15 CONSENSUS: EDENSKI-x-ELIVMETI-x(3)-ELIVMFTI(2)-x(b)-G-K-EKRI-

x(5)-ELIVMFI-ELIVMFCI-

CONSENSUS: x(2)-ESTAJ.

NAME: Ribosomal protein S& signature.

20 CONSENSUS: EGED-x(2)-ELIVD(2)-ESTYD-T-x(2)-G-ELIVMD(2)-x(4)-

EAGI-EKRHAYII.

NAME: Ribosomal protein S9 signature.

CVATZOJ - x(Z) - x(Z)

25 EKRI-EGSALI-ELIFI.

NAME: Ribosomal protein SLO signature.

CDNZENZUZ: CDNZRJ-ELIVMJ-x-

ELIVMI-P-T.

30

NAME: Ribosomal protein Sll signature.

CONSENSUS: ELIVMFD-x-EGSTACD-ELIVMFD-x(2)-EGSTALD-x(0,1)-

EGSNJ-ELIVMFJ-x-ELIVMJ-

CONSENSUS: x(4)-EDENJ-x-T-P-x-EPAJ-ESTCHJ-EDNJ.

35

NAME: Ribosomal protein Sl2 signature.

CONSENSUS: ERKI-x-P-N-S-EARI-x-R.

NAME: Ribosomal protein Sl3 signature.

40 CONSENSUS: EKRQSI-G-x-R-H-x(2)-EGSNHI-x(2)-ELIVMCI-R-G-Q.

NAME: Ribosomal protein S14 signature.

 $\mathbb{E}\mathsf{RGJ}-\mathsf{x}(3)-\mathbb{E}\mathsf{RNJ}$

45

NAME: Ribosomal protein S15 signature.

CONSENSUS: ELIVMI-x(2)-H-ELIVMFYI-x(5)-D-x(2)-ESAGNI-x(3)-

[LF]-x(9)-[LIVM]-x(2)-

CONSENSUS: EFYI.

50

NAME: Ribosomal protein SLL signature.

CONSENSUS: ELIVMID-X-ELIVMID-EKRID-L-ESTAKID-R-x-G-EAKRID.

NAME: Ribosomal protein S17 signature.

55 CONSENSUS: G-D-x-ELIVD-x-ELIVAD-x-EQEKD-x-ERKD-P-ELIVD-S.

NAME: Ribosomal protein SLB signature.

CONSENSUS: EIVJ-EYYJ-Y-X(2)-ELTYNTJ-X(2)-ELTYNTJ-X(2)-EFYTJ-

[LIVM]-[ST]-[DERP]-x-

CONSENSUS: EGYD-K-ELIVMD-x(3)-R-ELIVMASD.

5 NAME: Ribosomal protein S19 signature.

CONSENSUS: ESTANDID-G-EKRAND-x(b)-ELIVND-x(4)-ELIVND-EGSDD-

10 NAME: Ribosomal protein S21 signature.

CONSENSUS: EDED-x-A-ELYD-EKRD-R-F-K-EKRD-x(3)-EKRD.

NAME: Ribosomal protein S3Ae signature.

CONSENSUS: ELIVI-x-EGHI-R-EIVI-x-E-x-ESCI-L-x-D-L.

15

NAME: Ribosomal protein S4e signature.

CONSENSUS: H-x-K-R-CLIVMJ-CSANJ-x-P-x(2)-U-x-CLIVMJ-x-CKRJ.

NAME: Ribosomal protein She signature.

20 CONSENSUS: ELIVMI-ESTAMRIJ-G-G-x-D-x(2)-G-x-P-M.

NAME: Ribosomal protein S7e signature.

CONSENSUS: EKRJ-L-x-R-E-L-E-K-K-F-ESAPJ-x-EKRJ-H.

25 NAME: Ribosomal protein See signature.

CONSENSUS: R-x(2)-T-G-EGAI-x(5)-EHRI-K-EKRI-x-K-x-E-ELMI-G.

NAME: Ribosomal protein Sl2e signature.

-J-D-ENGI-(E)x-EAZI-(S)x-J-V-x-EQDRNI-LA

30

NAME: Ribosomal protein Sl7e signature.

CONSENSUS: A-x-I-x-ESTJ-K-x-L-R-N-EKRJ-I-A-G-EFYJ-x-T-H.

NAME: Ribosomal protein Sl9e signature.

35 CONZENSUS: P-x(b)-ECAN3-x(2)-ELIVMA3-x-R-x-EALIV3-ELV3-Q-x-L-

EEQI-

NAME: Ribosomal protein S2le signature.

CONZENSUS: L-Y-V-P-R-K-C-S-ESAJ.

40

NAME: Ribosomal protein S24e signature.

CONSENSUS: EFAD-G-x(2)-EKRJ-ESTAJ-x-G-EFYJ-EGAJ-x-ELIVMJ-Y-

- ENZJ-ENGJ

45 NAME: Ribosomal protein S2Le signature.

CONSENSUS: EYHI-C-V-S-C-A-I-H.

NAME: Ribosomal protein S27e signature.

CONZENZUZ: $\mathbb{L}QK\mathbb{J}-C-x(2)-C-x(b)-F-\mathbb{L}GZ\mathbb{J}-x-\mathbb{L}DZY-x(2)-C-x(2)-C-x(2)$

SO = CSD - x(2) - L - x(2) - P - x - G

NAME: Ribosomal protein S28e signature.

CONSENSUS: E-ESTI-E-R-E-A-R-x-L.

55 NAME: DNA mismatch repair proteins mutL / hexB / PMSL

signature.

CONSENSUS: G-F-R-G-E-A-L.

NAME: DNA mismatch repair proteins mutS family signature. CONSENSUS: ESTD-ELIVMD-x-ELIVMD-x-D-E-ELIVMD-EGCD-ERKHD-G-EGSTD-x(4)-G.

5 NAME: mutT domain signature.
CONSENSUS: G-x(5)-E-x(4)-ESTAGCI-ELIVMACI-x-R-E-ELIVMFTI-x-EE.

NAME: DnaA protein signature.

10 CONSENSUS: I-EGAI-x(2)-ELIVMFI-ESGNKI-x(0-1)-EKRI-x-H-ESTPI-ESTVI-ELIVMI(2)-xCONSENSUS: ESAI-x(2)-EKREI-ELIVMI.

NAME: Small, acid-soluble spore proteins, alpha/beta type, signature 1.
CONSENSUS: K-x-E-ELIVI-A-x-EDEI-ELIVMFI-G-ELIVMFI.

NAME: Small, acid-soluble spore proteins, alpha/beta type, signature 2.

20 CONSENSUS: EKRI-ESAQI-x-G-x-V-G-G-x-ELIVMI-x-EKRI(2)-ELIVMI(2).

NAME: Zinc-containing alcohol dehydrogenases signature-CONSENSUS: G-H-E-x(2)-G-x(5)-EGAJ-x(2)-EIVSACJ.

NAME: Quinone oxidoreductase / zeta-crystallin signature.

CONSENSUS: EGSDJ-EDEQHJ-x(2)-L-x(3)-ESAJ(2)-G-G-x-G-x(4)-Q-x(2)-EKRJ.

30 NAME: Iron-containing alcohol dehydrogenases signature 1.

CONSENSUS: ESTALIVI-ELIVFI-x-EDEI-x(6,7)-P-x(4)-EALIVI-x
EGSTI-x(2)-D-ETAIVMI
CONSENSUS: ELIVMFI-x(4)-E.

- 40 NAME: Short-chain dehydrogenases/reductases family signature.

 CONSENSUS: ELIVSPADNKI-x(12)-Y-EPSTAGNCVI-ESTAGNQCIVMI-ESTAGCI-K-{PC}-ESAGFRI-CONSENSUS: ELIVMSTAGDI-x(2)-ELIVMFYWI-x(3)-ELIVMFYWGAPTHQI-45 EGSACQRHMI.

NAME: Aldo/keto reductase family signature L.
CONSENSUS: G-EFYI-R-EHSALI-ELIVMFI-D-ESTAGCI-EASI-x(5)-Ex(2)-ELIVMI-G.

50
NAME: Aldo/keto reductase family signature 2.
CONSENSUS: ELIVMFY3-x(9)-EKREQ3-x-ELIVM3-E-ELIVM3-ESC3-NEFY3.

55 NAME: Aldo/keto reductase family putative active site signature.

CONSENSUS: CLIVMD-CPAIVD-CKRD-CSTD-x(4)-R-x(2)-CGSTAEQKD-CNSLD-x(2)-CLIVMFAD.

Homoserine dehydrogenase signature.

NAME:

 $-(5)\times -(5)\times -(5)$ **CONZENZUZ:** CLIVMB-x-G-x-D-x(3)-K. 5 NAME: NAD-dependent glycerol-3-phosphate dehydrogenase signature. G-EATD-ELIVMD-K-EDND-ELIVMD(2)-A-x-EGAD-x-G-**CONSENSUS:** CLIVMFD-x-CDED-G-CLIVMD-x-**CONSENSUS: ELIVMFYWD-G-x-N**. 10 NAME: FAD-dependent glycerol-3-phosphate dehydrogenase signature 1. **CONSENSUS:** 15 NAMF: FAD-dependent glycerol-3-phosphate dehydrogenase signature 2. **CONZENZUZ:** G-G-K-x(2)-EGSTEJ-Y-R-x(2)-A-Mannitol dehydrogenases signature. 20 NAME: ELIVMYD-x-EFSD-x(2)-ESTAGCVD-x-V-D-R-EIVD-x-EPSD. **CONSENSUS:** Histidinol dehydrogenase signature. NAME: **CONSENSUS:** I-D-x(2)-A-G-P-EZTJ-E-ELIVJJ-ELCVIJ-B-G-D-A-G-D-X(3)-25 A-x(4)-ELIVMI-EAVI-ESACLD-EDED-ELIVMFCD-ELIVMD-ESAD-x(2)-E-H. CONZENZUZ: L-lactate dehydrogenase active site. NAME: CONSENSUS: ELIVMAD-G-EEQD-H-G-EDND-ESTD. 30 D-isomer specific 2-hydroxyacid dehydrogenases NAD-NAME: binding signature. **CONSENSUS:** CLIVMA3-CAG3-CIVT3-CLIVMFY3-EAG3-x-G-ENHKRQGSAC3-ELIV3-G-x(13-14)-35 **CONZENZUZ:** ELIVINT3-x(2)-EFYwCTH3-EDNSTK3. NAME: D-isomer specific 2-hydroxyacid dehydrogenases signature 2. CONSENSUS: ELIVMFYWAD-ELIVFYWCD-x(2)-ESACD-EDNQHRD-EIVFAD-40 CLIVFD-x-CLIVFD-CHNID-x-P-x(4)-ESTND-x(2)-ELIVMFD-x-EGSDND: CONZENZUZ: D-isomer specific 2-hydroxyacid dehydrogenases NAME: signature 3. ELMFATCJ-EKPQJ-x-EGSTDNJ-x-ELIVMFYWRJ-CONSENSUS: ELIVMFYUJ(2)-N-x-ESTAGCJ-R-EGPJ-x-**CONSENSUS:** ELIVHI-ELIVMCI-EDNVI. 3-hydroxyisobutyrate dehydrogenase signature. NAMF: TLIVMFYD(2)-G-L-G-x-EMQD-G-x-EPGSD-EMAD-TSAD. 50 **CONZENZUZ:** Hydroxymethylglutaryl-coenzyme A reductases signature NAME: l. **CONSENSUS:** [RKH]-x(b)-D-x-M-G-x-N-x-[LIVMA]. 55 Hydroxymethylglutaryl-coenzyme A reductases signature NAME: 2. ELIVMD-G-x-ELIVMD-G-G-EAGD-T. **CONSENSUS:**

NAME: Hydroxymethylglutaryl-coenzyme A reductases signature 3.

CONSENSUS: A-ELIVMI-x-ESTANI-x(2)-ELII-x-EKRNQI-EGSAI-H-ELMI-

5 x-EFYLHD.

NAME: Hydroxymethylglutaryl-coenzyme A reductases profile.

NAME: 3-hydroxyacyl-CoA dehydrogenase signature.

10 CONZENZUS: EDNEI-x(2)-EGAI-F-ELIVMFYI-x-ENTI-R-x(3)-EPAI-ELIVMFYI(2)-x(5)-

CONSENSUS: ELIVMFYCTD-ELIVMFYD-x(2)-EGVD.

NAME: Malate dehydrogenase active site signature.

15 CONSENSUS: ELIVND-T-EMNXATD-T-EMVID-L-00-x(2)-R-ESTAD-x(2)-ELIVND-T-EMVID-

NAME: Malic enzymes signature.

CONSENSUS: F-x-EDVJ-D-x(2)-G-T-EGSAJ-x-ELVJ-x-ELIVMAJ-

EGASTI(2)-ELIVMFI(2).

NAME: Isocitrate and isopropylmalate dehydrogenases signature.

-CNDST-Z-X-UZVMID-CTNDD-B-CNDVD-ELTYMID-CZNZVD-X-EZNDVD-CNDZ-

-0-(4,E)x-E9A2J-(5)x

25 CONSENSUS: ESTGI-ELIVMPAI-G-ELIVMFI.

NAME: b-phosphogluconate dehydrogenase signature.
CONSENSUS: ELIVAD-x-0-x(2)-EGDD-K-G-T-G-x-W.

30 NAME: Glucose-b-phosphate dehydrogenase active site-CONSENSUS: D-H-Y-L-G-K-TEQKI.

NAME: IMP dehydrogenase / GMP reductase signatureCONSENSUS: ELIVMI-ERKI-ELIVMI-G-ELIVMI-G-x-G-S-ELIVMI-C-x-T.

NAME: Bacterial quinoprotein dehydrogenases signature l. CONSENSUS: EDENI-W-x(3)-G-ERKI-x(6)-EFYWI-S-x(4)-ELIVMI-N-x(2)-N-V-x(2)-L-ERKI.

40 NAME: Bacterial quinoprotein dehydrogenases signature 2. CONSENSUS: W-x(4)-Y-D-x(3)-EDNJ-ELIVMFYJ(4)-x(2)-G-x(2)-ESTAJ-P.

NAME: FMN-dependent alpha-hydroxy acid dehydrogenases active

45 site. CONSENSUS: S-N-H-G-EAGD-R-Q.

NAME: GMC oxidoreductases signature 1.

CENSUS: EQUIDENCE CONSUME CONSUME CAND CONSUME CAND CONSUMERADIO CONSUMERA C

50 EFYWAJ-x(2)-EPAGJ-x(5)-CONSENSUS: EDNESHJ-

NAME: GMC oxidoreductases signature 2.

55 ELIVMI-G.

NAME: Eukaryotic molybdopterin oxidoreductases signature.

CONSENSUS: EGAJ-x(3)-ECTHQNRXXJ-x(A)-ELIVARYWSJ-x(A)-

ELIVMFI-x-C-x(2)-EDENI-RCONSENSUS: x(2)-EDEII-

5 NAME: Prokaryotic molybdopterin oxidoreductases signature l. CONSENSUS: ESTANJ-x-ECHJ-x(2,3)-C-ESTAGJ-EGSTVMFJ-x-C-x-

CLIVMFYWJ-x-ELIVMAJ-x(3,4)CONSENSUS:

NAME: Prokaryotic molybdopterin oxidoreductases signature 2-CONSENSUS: ESTAI-x-ESTACI(2)-x(2)-ESTAI-D-ELIVMYI(2)-L-P-x-ESTACI(2)-x(2)-E.

NAME: Prokaryotic molybdopterin oxidoreductases signature 3.

15 CONSENSUS: A-x(3)-EGDTJ-I-x-EDRQTKJ-x-EDEAJ-x-ELIVMJ-x-ELIVMJ-x-ELIVMJ-x-EDRJ-x-ENSJ-x(2)-EGSJ-

CONSENSUS: x(5)-A-x-ELIVMI-ESTI.

NAME: Aldehyde dehydrogenases glutamic acid active site.

20 CONSENSUS: ELIVMFGAJ-E-ELIMSTACJ-EGSJ-G-EKNLMJ-ESADNJETAPFVJ.

25 EGSTADNEKRI.

NAME: Aspartate-semialdehyde dehydrogenase signatureCONSENSUS: ELIVMI-ESADNI-x(2)-C-x-R-ELIVMI-x(4)-EGSCI-HESTAI.

ESTAI 30

NAME: Glyceraldehyde 3-phosphate dehydrogenase active site.
CONSENSUS: EASUB-S-C-ENTB-T-x(2)-ELIMB.

NAME: N-acetyl-gamma-glutamyl-phosphate reductase active site.

35 site. CONSENSUS: ELIVMI-EGSAI-x-P-G-C-EFYI-EAVPI-T-EGAI-x(3)-EGYAI-EUVII-x-P-

NAME: Gamma-glutamyl phosphate reductase signature.

40 CONSENSUS: V-x(5)-A-ELIV3-x-H-I-x(2)-EHY3-EGS3-EST3-x-H-EST3-

EDE3-x-I.

NAME: Dihydrodipicolinate reductase signature.
CONSENSUS: E-EIVI-x-E-x-H-x(3)-K-x-D-x-P-S-G-T-A.

45

50 NAME: Dihydroorotate dehydrogenase signature 2.

CONSENSUS: ELIVI(2)-EGSAI-x-G-G-EIVI-x-ESTGNI-x(3)-EACVIx(b)-G-A.

NAME: Coproporphyrinogen III oxidase signature.

55 CONSENSUS: K-x-W-C-x(2)-EFYHI(3)-ELIVMI-x-H-R-x-E-x-R-G-ELIVMI-G-G-ELIVMI-F-F-D.

NAME: Fumarate reductase / succinate dehydrogenase FAD-

binding site.

CONSENSUS: R-ESTJ-H-ESTJ-x(2)-A-x-G-G.

5 NAME: Acyl-CoA dehydrogenases signature 1.

EGSAI.

NAME: Acyl-CoA dehydrogenases signature 2.

10 CONSENSUS: EQDED-x(2)-G-EGSD-x-G-ELIVMFYD-x(2)-EDEND-x(4)-

EKRI-x(3)-EDENI.

NAME: Alanine dehydrogenase & pyridine nucleotide

transhydrogenase signature 1.

15 CONSENSUS: G-ELTVMI-P-x-E-x(3)-N-E-x(1-3)-R-V-A-x-ESTI-P-x-

LCSTI-A-x(5)-F-x-EKHI-

CONZENZUZ: x-G.

NAME: Alanine dehydrogenase & pyridine nucleotide

20 transhydrogenase signature 2.

-x-CADD-(E)x-CAZD-(S)x-D-A-x-D-EADD-D-(2)EMVID

-V-x-A-Q-EMVIJI-EQZJ -d-(E)x :ZUZNJZNO)

25 NAME: Glu / Leu / Phe / Val dehydrogenases active site-

CONSENSUS: ELIVI-x(2)-G-G-ESAGI-K-x-EGVI-x(3)-EDNSTI-EPLI.

NAME: D-amino acid oxidases signature.

-CUZNAZNO: ELIVIDI(2)-H-ENHAD-Y-G-x-EQADI(2)-x-G-x(5)-G-x-A-

30

NAME: Pyridoxamine 5'-phosphate oxidase signature.

CONSENSUS: ELIVFI-E-F-W-EQHGI-x(4)-R-ELIVMI-H-EDNEI-R.

NAME: Copper amine oxidase topaquinone signature.

35 CONSENSUS: ELIVIDATA ELIVIDATA (4)-T-x(2)-N-Y-EDEJ-EVIJ.

NAME: Copper amine oxidase copper-binding site signature.

CONSENSUS: T-x-G-x(2)-H-ELIVMFJ-x(3)-E-EDEJ-x-P.

40 NAME: Lysyl oxidase putative copper-binding region

signature.

CONSENSUS: W-E-W-H-S-C-H-Q-H-Y-H.

NAME: Delta 1-pyrroline-5-carboxylate reductase signature.

45 CONSENSUS: EPALFI-x(2,3)-ELVIJ-E(E)x-E71AQJ-ESTVI-x-

EGANB-G-x-T-x(2)-EAGD-

CONSENSUS: ELIVI-x(2)-ELMFI-EDENGKI.

NAME: Dihydrofolate reductase signature.

50 CONSENSUS: ELVAGCI-ELIFI-G-x(4)-ELIVMFI-P-W-x(4,5)-EDEI-x(3)-

EFYIVI-x(3)-ESTIQI.

NAME: Tetrahydrofolate dehydrogenase/cyclohydrolase

signature l.

55 CONSENSUS: EEQI-x-EEQKI-ELIVMI(2)-x(2)-ELIVMI-x(2)-ELIVMYI-N-

x-EDNJ-x(5)-ELIVMFJ(3)-

CONSENSUS: Q-L-P-ELVI.

NAME: Tetrahydrofolate dehydrogenase/cyclohydrolase

signature 2.

CONSENSUS: P-G-G-V-G-P-EMFI-T-EIVI.

5 NAME: Oxygen oxidoreductases covalent FAD-binding site.
CONSENSUS: P-x(10)-EDEJ-ELIVMJ-x(3)-ELIVMJ-x(9)-ELIVMJ-x(3)EGSAJ-EGSTJ-G-H.

NAME: Pyridine nucleotide-disulphide oxidoreductases class-I

10 active site.

25

55

CONSENSUS: G-G-x-C-ELIVAI-x(2)-G-C-ELIVMI-P.

NAME: Pyridine nucleotide-disulphide oxidoreductases class-II active site.

15 CONSENSUS: $C-x(2)-C-D-\mathbb{E}GA\mathbb{I}-x(2-4)-\mathbb{E}FY\mathbb{I}-x(4)-\mathbb{E}LIVM\mathbb{I}-x-\mathbb{E}LIVM\mathbb{I}(2)-G(3)-\mathbb{E}DN\mathbb{I}$.

NAME: Respiratory-chain NADH dehydrogenase subunit 1 signature 1.

20 CONSENSUS: G-ELIVMFYKSJ-ELIVMGPJ-Q-x-ELIVMFYJ-x-D-EAGIMJ-ELIVMFYJ-K-ELVMYSTJ-CONSENSUS: ELIVMFYGJ-x-EKRJ-EEQGJ.

NAME: Respiratory-chain NADH dehydrogenase subunit lesignature 2.

CONSENSUS: P-F-D-ELIVMFYQD-ESTAGPVMD-E-EQD-ELVMSD-x(2)-G.

NAME: Respiratory-chain NADH dehydrogenase 20 Kd subunit) signature.

30 signature.
CONSENSUS: EGNI-x-D-EKRSTI-ELIVMFI(2)-P-EIVI-D-ELIVMFYWI(2)x-P-x-C-P-EPTI.

NAME: Respiratory-chain NADH dehydrogenase 24 Kd subunit signature.
CONSENSUS: D-x(2)-F-EST3-x(5)-C-L-G-x-C-x(2)-EGA3-P.

NAME: Respiratory chain NADH dehydrogenase 3D Kd subunit signature.

40 CONSENSUS: E-R-E-x(2)-EDED-ELIVMF3(2)-x(6)-EHKD-x(3)-EKRPD-x-ELIVMD-ELIVMS3.

NAME: Respiratory chain NADH dehydrogenase 49 Kd subunit signature.

45 CONSENSUS: ELIVMHI-H-ERTI-EGAI-x-E-K-ELIVMTI-x-E-x-EKRQI.

NAME: Respiratory-chain NADH dehydrogenase 51 Kd subunit signature 1.
CONSENSUS: G-EAMI-G-EARI-Y-ELIVMI-C-G-EDEI(2)-ESTAI(2)-

50 ELIMD(2)-EEND-S.

NAME: Respiratory-chain NADH dehydrogenase 51 Kd subunit signature 2.
CONSENSUS: E-S-C-G-x-C-x-P-C-R-x-G.

NAME: Respiratory-chain NADH dehydrogenase 75 Kd subunit signature LCONSENSUS: P-x(2)-C-EYWSI-x(7)-G-x-C-R-x-C

NAME: Respiratory-chain NADH dehydrogenase 75 Kd subunit signature 2.

CONSENSUS: C-P-x-C-DEJ-x-EGSJ(2)-x-C-x-L-Q.

NAME: Respiratory-chain NADH dehydrogenase 75 Kd subunit signature 3.
CONSENSUS: R-C-ELIVMD-x-C-x-R-C-ELIVMD-x-EFYD.

10 NAME: Nitrite and sulfite reductases iron-sulfur/siroheme-binding site.

CONSENSUS: ESTVI-G-C-x(3)-C-x(b)-EDEI-ELIVMFI-EGATI-ELIVMFI.

NAME: Uricase signature.

15 CONSENSUS: L-x-ELV3-L-K-EST3-T-x-S-x-F-x(2)-EFY3-x(4)-EFY3.

NAME: Heme-copper oxidase catalytic subunit, copper B binding region signature.

CONSENSUS: EYWGI-ELIVFYWTAI(2)-EVGSI-H-ELNPI-x-V-x(44,47)-H-

20 H-

NAME: CO II and nitrous oxide reductase dinuclear copper centers signature.

CONSENSUS: V-x-H-x(33,40)-C-x(3)-C-x(3)-H-x(2)-M.

NAME: Cytochrome c oxidase subunit Vb, zinc binding region signature.

CONSENSUS: $\mathbb{L}IVMJ(2) - \mathbb{L}FYMJ - x(JD) - C - x(2) - C - G - x(2) - \mathbb{L}FYJ - K - L$.

30 NAME: Multicopper oxidases signature 1.

CONSENSUS: G-x-EFYWI-x-ELIVMFYWI-x-ECSTI-x(8)-G-ELMI-x(3)
ELIVMFYWI.

NAME: Multicopper oxidases signature 2.
35 CONSENSUS: H-C-H-x(3)-H-x(3)-EAGI-ELMI.

NAME: Peroxidases proximal heme-ligand signature.

CONSENSUS: EDETD-ELIVMID-x(2)-ELIVMID-ELIVMSTAGD-ESAGDELIVMSTAGD-H-ESTAD-ELIVMYD.

NAME: Peroxidases active site signature.

CONSENSUS: ESGATUB-x(3)-ELIVMAB-R-ELIVMAB-x-EFWB-H-x-ESACB.

NAME: Catalase proximal heme-ligand signature.
45 CONSENSUS: R-ELIVMFSTANI-F-EGASTNPI-Y-x-D-EASTI-EQEHI.

NAME: Glutathione peroxidases selenocysteine active site.
CONSENSUS: [[GN]-[RKHNFYC]-x-[LIVMFC]-[LIVMF](2)-x-N-[VT]-x[[STC]-x-C-[[GA]-x-T.

55 NAME: Glutathione peroxidases signature 2. CONSENSUS: ELIVD-EAGDD-F-P-ECSD-ENGD-Q-F.

NAME: Lipoxygenases iron-binding region signature 1.

CONSENSUS: H-EEQ3-x(3)-H-x-ELM3-ENQRC3-EGST3-H-ELIVMSTAC3(3)-

NAME: Extradiol ring-cleavage dioxygenases signature.

CONSENSUS: EGNTIVI-x-H-x(5,7)-ELIVMFI-Y-x(2)-EDENTAI-P-x-

IGPI-x(2,3)-E.

15

NAME: Indoleamine 2-3-dioxygenase signature l-CONSENSUS: G-G-S-EANJ-EGAJ-Q-S-S-x(2)-Q.

NAME: Indoleamine 2,3-dioxygenase signature 2.

20 CONSENSUS: EFY3-L-EDQ3-EDE3-ELIVM3-x(2)-Y-M-x(3)-H-EKR3.

NAME: Bacterial ring hydroxylating dioxygenases—alphasubunit signature.

CONSENSUS: C-x-H-R-EGAI-x(8)-G-N-x(5)-C-x-EFYI-H.

25

NAME: Bacterial luciferase subunits signature.

CONSENSUS: EGAB-ELIVMB-P-ELIVMB-x-ELIVMFYB-x-W-x(b)-ERKB-x(b)-Y-x(3)-EARB.

30 NAME: ubiH/COQL monooxygenase family signature-CONSENSUS: H-P-ELIVI-EAGI-G-Q-G-x-N-x-G-x(2)-D.

NAME: Biopterin-dependent aromatic amino acid hydroxylases signature.

35 CONSENSUS: P-D-x(2)-H-EDEJ-ELIJ-ELIVMFJ-G-H-ELIVMCJ-P.

NAME: Copper type II, ascorbate-dependent monooxygenases signature 1.

CONSENSUS: H-H-M-x(2)-F-x-C.

40

NAME: Copper type II, ascorbate-dependent monooxygenases signature 2.
CONSENSUS: H-x-F-x(4)-H-T-H-x(2)-G.

45 NAME: Tyrosinase CuA-binding region signature.
CONSENSUS: H-x(4,5)-F-ELIVMFTP3-x-EFW3-H-R-x(2)-ELM3-x(3)-E.

NAME: Tyrosinase and hemocyanins CuB-binding region signature.

.d-(E)x-H-(E)x-EUYTMVIJ-7-x-9-0

NAME: Fatty acid desaturases family 1 signature. CONSENSUS: G-E-x-EFY1-H-N-EFY1-H-H-x-F-P-x-D-Y.

NAME: Cytochrome P450 cysteine heme-iron ligand signature.
CONSENSUS: EFWI-ESGNHI-x-EGDI-x-ERHPTI-x-C-ELIVMFAPI-EGADI.

NAME: Heme oxygenase signature.
5 CONSENSUS: L-L-V-A-H-A-Y-T-R.

20

NAME: Copper/Zinc superoxide dismutase signature 1.
CONSENSUS: EGAI-EIFATI-H-ELIVFI-H-x(2)-EGPI-ESDGI-x-ESTAGDI.

10 NAME: Copper/Zinc superoxide dismutase signature 2. CONSENSUS: G-EGNI-ESGAI-G-x-R-x-ESGAI-C-x(2)-EIVI.

NAME: Manganese and iron superoxide dismutases signature.
CONSENSUS: D-x-W-E-H-ESTAI-EFYI(2).

NAME: Ribonucleotide reductase large subunit signature.
CONSENSUS: W-x(2)-ELFJ-x(6,7)-G-ELIVMJ-EFYRAJ-ENHJ-x(3)ESTAGLIVMJ-EASCJ-x(2)CONSENSUS: EPAJ.

NAME: Ribonucleotide reductase small subunit signature.
CONSENSUS: EIVMSEQI-E-x(l,2)-ELIVTAI-EHYI-EGSAI-x-ESTAVMI-yx(2)-ELIVMQI-x(3)CONSENSUS: ELIFYI-EIVFYCSAI.

NAME: Nitrogenases component l alpha and beta subunits signature l.

CONSENSUS: ELIVMFYHI-ELIVMFSTI-H-EAGI-EAGSPI-ELIVMQAI-EAGI-C.

NAME: Nitrogenases component l alpha and beta subunits signature 2.

CONSENSUS: ESTANQI-EETI-C-x(5)-G-D-EDNI-ELIVMTI-x-ESTAGRI-ELIVMFYSTI.

NAME: NifH/frxC family signature 1.
CONSENSUS: E-x-G-G-P-x(2)-EGAJ-x-G-C-EAGJ-G.

NAME: NifH/frxC family signature 2.
40 CONSENSUS: D-x-L-G-D-V-V-C-G-G-F-EAGD-x-P.

NAME: Nickel-dependent hydrogenases large subunit signature 1.
CONSENSUS: R-G-ELIVMFI-E-x(15)-EQESMI-R-x-C-G-ELIVMI-C.

45
NAME: Nickel-dependent hydrogenases large subunit signature 2.
CONSENSUS: [FY]-D-P-C-[LIM]-[ASG]-C-x(2,3)-H.

50 NAME: Glutamyl-tRNA reductase signature.

CONSENSUS: H-ELIVMI-x(2)-ELIVMI-EGSTACI(3)-ELIVMI-EDEQI-S
ELIVMAI-ELIVMI(2)-EGFI-E
CONSENSUS: x-EQRI-EIVI-ELITI-ESTAGI-Q-ELIVMI-EKRI.

NAME: Bacterial-type phytoene dehydrogenase signature.

CONSENSUS: ENGI-x-EFYWVI-ELIVMFI-x-G-EAGCI-EGSI-ETAI-EHQTI-PG-ESTAVI-G-ELIVMICONSENSUS: x(5)-EGSI.

NAME: Glycine radical signature.

20

CONSENSUS: ESTIVI-x-R-EIVTI-ECSAI-G-Y-x-EGACVI.

5 NAME: Ergosterol biosynthesis ERG4/ERG24 family signature 1. CONSENSUS: G-x(2)-ELIVMD-Y-D-x-EFYD-x-G-x(2)-L-N-P-R.

NAME: Ergosterol biosynthesis ERG4/ERG24 family signature 2. CONSENSUS: £LIVMI(2)-H-R-x(2)-R-D-x(3)-C-x(2)-K-Y-G.

10 CONZENZOZ: FLIAMM(S)-H-K-X(S)-K-D-X(3)-C-X(S)-K-A-Q.

NAME: NNMT/PNMT/TEMT family of methyltransferases signature.
CONSENSUS: L-I-D-I-G-S-G-P-T-EIVB-Y-Q-L-L-S-A-C.

NAME: RNA methyltransferase trmA family signature 1.

15 CONSENSUS: EDNI-P-EPAI-R-x-G-x(14,16)-ELIVMI(2)-Y-x-S-C-N-x(2)-T.

NAME: RNA methyltransferase trmA family signature 2. CONSENSUS: ELIVMFI-D-x-F-P-EQHYI-ESTI-x-H-ELIVMFYI-E-

NAME: Thymidylate synthase active site-CONSENSUS: R-x(2)-ELIVMJ-x(3)-EFWJ-EQNJ-x(8,9)-ELVJ-x-P-C-EHAVMJ-x(3)-EQMTJ-EFYWJ-CONSENSUS: x-ELVJ-

NAME: Ribosomal RNA adenine dimethylases signature.

CONSENSUS: ELIVMI-ELIVMYI-EDEI-x-G-ESTAPVI-G-x-EGAI-xELIVMFI-ESTI-x(2)-ELIVMICONSENSUS: x(b)-ELIVMYI-x-ESTAGVI-ELIVMFYHCI-E-x-D.

NAME: Methylated-DNA--protein-cysteine methyltransferase active site.
CONSENSUS: [LIVMF]-P-C-H-R-[LIVMF](2).

35 NAME: N-b Adenine-specific DNA methylases signature. CONSENSUS: ELIVMACI-ELIVFYWAI-x-EDNI-P-P-EFYWI.

NAME: N-4 cytosine-specific DNA methylases signature. CONSENSUS: ELIVMFI-T-S-P-P-EFYI.

AO

NAME: C-5 cytosine-specific DNA methylases active site.

CONSENSUS: EDENKSD-x-EFLIVD-x(2)-EGSTCD-x-P-C-x(2)-EFYWLIMDS.

45 NAME: C-5 cytosine-specific DNA methylases C-terminal signature.

CONSENSUS:

CRKQGTF3-x(2)-G-N-ESTAG3-ELIVMF3-x(3)-ELIVMT3-x(3)-ELIVM3.

NAME: Uroporphyrin-III C-methyltransferase signature 2.

x(5,6)-ELIVMFYWPACD-

CONSENSUS: x-ELIVMYI-x-P-G.

5 NAME: ubiE/COQ5 methyltransferase family signature l-CONSENSUS: Y-D-x-M-N-x(2)-ELIVMI-S-x(3)-H-x(2)-W-

NAME: ubiE/COQ5 methyltransferase family signature 2.
CONSENSUS: R-V-ELIVMD-K-EPVD-G-G-x-ELIVMFD-x(2)-ELIVMD-E-x-S.

NAME: Serine hydroxymethyltransferase pyridoxal-phosphate attachment site.

CONSENSUS: EDEHJ-ELIVMFYJ-x-ESTMVJ-EGSTJ-ESTJ(2)-H-K-ESTJ-

NAME: Phosphoribosylglycinamide formyltransferase active site.

CONSENSUS: G-x-ESTMD-EIVTD-x-EFYWVQD-EVMD-x-EDEVMD-x-

20 ELIVMYI-D-x-G-x(2)-ELIVTI-CONSENSUS: x(b)-ELIVMI.

NAME: Aspartate and ornithine carbamoyltransferases signature.

25 CONSENSUS: F-x-EEKJ-x-S-EGTJ-R-T.

NAME: Transketolase signature 1.

CONSENSUS: R-x(3)-ELIVMTAJ-EDENQSTHKFJ-x(5,6)-EGSNJ-G-H-EPLIVMFJ-EGSTAJ-x(2)-

30 CONSENSUS: ELIMCD-EGSD.

NAME: Transketolase signature 2.

CONSENSUS: G-EDEQGSAJ-EDNJ-G-EPAEQJ-ESTJ-EHQJ-x-EPAGMJELIVMYACJ-EDEFYWJ-x(2)-

35 CONSENSUS: ESTAPD-x(2)-ERGAD.

NAME: Transaldolase signature l.

CONSENSUS: EDGI-EIVSAI-T-ESTI-N-P-ESTAI-ELIVMFI(2).

40 NAME: Transaldolase active siteCONSENSUS: ELIVMI-x-ELIVMI-K-ELIVMI-EPASI-x-ESTI-x-EDENGPASIG-ELIVMI-x-EAGVI-xCONSENSUS: EGEKRSTI-x-ELIVMI-

45 NAME: Acyltransferases ChoActase / COT / CPT family signature 1.

CONSENSUS: ELID-P-x-ELVPD-P-EIVTAD-P-x-ELIVMD-x-EDENGASD-ESTD-ELIVMD-x(2)-ELYD.

50 NAME: Acyltransferases ChoActase / COT / CPT family signature 2.

CONSENSUS: R-EFYWJ-x-EDAJ-EKAJ-x(0,1)-ELIVMFYJ-x-ELIVMFYJ(2)-x(3)-EDNSJ-EGSAJ-x(b)CONSENSUS: EDEJ-EHSJ-x(3)-EDEJ-EGAJ.

55
NAME: Thiolases acyl-enzyme intermediate signature.
CONSENSUS: ELIVMD-ENSTD-x(2)-C-ESAGLID-ESTD-ESAGD-ELIVMFYNSDx-ESAGD-ELIVMD-x(6)-

ELIVMI. **CONZENZUZ:**

Thiolases signature 2.

CONZENZUZ: -9-KTZJ-x-G-H-P-x-G-X-EKVIJJ-x-G-X-ETJ-G-

5

NAME: Thiolases active site.

CONSENSUS: EAGD-ELIVMAD-ESTAGLIVMD-ESTAGD-ELIVMAD-C-x-EAGD-x-

EAGD-x-EAGD-x-ESAGD.

10 NAME: Chloramphenicol acetyltransferase active site.

CONSENSUS: Q-ELIVJ-H-H-ESAJ-x(2)-D-G-EFYJ-H.

NAME: Hexapeptide-repeat containing-transferases signature. **CONSENSUS:** ELIVI-EGAEDI-x(2)-EVIJI-x-ELIVI-x(3)-ELIVACI-x-

15 ELIVI-EGAEDII-x(2)-

> CONSENSUS: CSTAVRD-x-CLIVD-CGAEDD-x(2)-CVATZD-(2)-

ELIV3.

NAME: Beta-ketoacyl synthases active site.

20 CONSENSUS: G-x(4)-ELIVMFAPD-x(2)-EAGCD-C-ESTAD(2)-ESTAGD-

x(3)-ELIVMFI.

NAMF: Chalcone and stilbene synthases active site.

R-ELIVMFYSD-x-ELIVMD-x-EQHGD-x-G-C-EFYNAD-EGAD-G-CONZENZUZ:

CGAJ-CSTAVJ-x-CLIVMFJ-25

CONZENZUZ:

NAME: Myristoyl-CoA: protein N-myristoyltransferase signature

30 **CONZENZUZ:** E-I-N-F-L-C-x-H-K.

> NAME: Myristoyl-CoA:protein N-myristoyltransferase signature

CONSENSUS: K-F-G-x-G-D-G.

35

NAME: Gamma-glutamyltranspeptidase signature.

CONZENZUS: -x-EAT23-V-x-EN23-D-(+)x-EAMVIJ3-ET23-x-H-CAT23-T

T-x-T-ELIVMD-ENED-

x(1,2)-EFY1-G. CONSENSUS:

40

Transglutaminases active site.

CONZENSUS: EGTJ-Q-CCAJ-W-V-x-CSAJ-CGAJ-CIVTJ-x(2)-T-x-CLMSCJ-

R-ECSAD-ELVD-G-

Phosphorylase pyridoxal-phosphate attachment site. 45

CONSENSUS: E-A-ESCI-G-x-EGSI-x-M-K-x(2)-ELMI-N.

NAME: UDP-glycosyltransferases signature.

CONZENZUZ: EFWJ-x(2)-Q-x(2)-ELIVMYAJ-ELIMVJ-x(4,6)-ELVGACJ-

50 ELVFYAD-ELIVMFD-ESTAGCMD-

> CONZENSUS: __EATMVIJ-EJDATZJ-(E)x-EDATZJ-(E)x-D-EJDATZJ-EDNHJ-

x(4)-EPQRI-ELIVMTI-

CONZENSUS: · CHHAD-CZ3CD-(E)x-CAGD-(E)x

55 NAME: Purine/pyrimidine phosphoribosyl transferases

signature.

CONSENSUS: ELIVMFYWCTAD-ELIVMD-ELIVMAD-ELIVMFCD-EDED-D-

-EGVAT23-ELWVIJ-EZMVIJ-

CONZENZUZ:

ESTARI-EGACI-x-ESTARI.

Glutamine amidotransferases class-I active site.

CONSENSUS: EPASD-ELIVMFYTD-ELIVMFYD-G-ELIVMFYD-C-ELIVMFYND-G-

x-EQEHD-x-ELIVMFAD.

Glutamine amidotransferases class-II active site. <x(O,11)-C-EGSI-EIVI-ELIVMFYWI-EAGI.</pre> CONSENSUS:

10 NAME: Purine and other phosphorylases family 1 signature. : CONZENZUZ EGTD-x-G-ELIVMD-G-x-EAT2D-x-2-EAT2D-x-G-ENVID-E-L.

Purine and other phosphorylases family 2 signature. CONSENSUS: $\mathbb{CLIVJ} - x(3) - G - x(2) - H - x - \mathbb{CLIVMFYJ} - x(4) - \mathbb{CLIVMFJ} - x(3) - \mathbb{CLIVMFJ}$

15 EATVD-x(1,2)-ELIVMD-x-CONZENZUZ: -CAZJ-CVTSJ-(4)-CVTSJ-(4)-CVTAJ-(4)-CVTAJx-G-EGSJ-ELIVMJ.

NAME: Thymidine and pyrimidine-nucleoside phosphorylases

20 signature.

30

45

CONZENZUZ: S-EGSJ-R-EGAD-ELIVD-x(2)-ETAD-EGAD-G-T-x-D-x-CLIVI-E.

NAME: ATP phosphoribosyltransferase signature.

25 **CONZENZUZ:** E-x(5)-G-x-ESAGJ-x(2)-EVJ-x-D-ELIVJ-x(2)-ESJJ-Gx-T-ELM3.

NAME: NAD:arginine ADP-ribosyltransferases signature. **CONZENZUZ:** EFYD-x-EFYD-K-x(2)-H-EFYD-x-L-ESTD-x-A.

Prolipoprotein diacylglyceryl transferase signature. **CONZENZUZ:** G-R-x-EGAB-N-F-ELIVMFD-N-x-E-x(2)-G.

NAME: S-adenosylmethionine synthetase signature 1. 35 CONZENZUZ: $G-A-G-D-Q-G-\times(3)-G-Y$.

NAMF: S-adenosylmethionine synthetase signature 2. CONZENZUZ: G-EGAI-G-EASCI-F-S-x-K-EDEI-

40 NAME: Polyprenyl synthetases signature 1. CONSENSUS: [LIVM](2)-x-D-D-x(2,4)-D-x(4)-R-R-[GH].

NAME: Polyprenyl synthetases signature 2. CONSENSUS: CLIVMFY3-G-x(2)-CFYL3-Q-CLIVM3-x-D-D-CLIVMFY3-x-EDNGI.

Squalene and phytoene synthases signature 1. CONSENSUS: Y-ECSAMI-x(2)-EVSGI-A-EGSAI-ELIVATI-EVID-G-x(2)-ELMZCI-x(2)-ELIVI.

50 Squalene and phytoene synthases signature 2. NAME: CONSENSUS: $\mathbb{L}IVM\mathbf{3}-G-\mathbf{x}(\mathbf{3})-Q-\mathbf{x}(\mathbf{2},\mathbf{3})-\mathbf{N}-\mathbb{L}IF\mathbf{3}-\mathbf{x}-\mathbf{R}-\mathbf{D}-\mathbb{L}IVMFY\mathbf{3}-\mathbf{x}(\mathbf{2})-$ EDED-x(4,7)-R-x-EFYD-CONZENZUZ: x-P.

55 NAME: Protein prenyltransferases alpha subunit repeat signature.

CONSCISS EPSIAVD-x-EVACADEVD-EVEQIYD-x-ELIVMAGPD-W-ENGSTHFD-EFYHQD-ELIVMRD.

NAME: Riboflavin synthase alpha chain family signature.

5 CONSENSUS: ELIVMF3-x(5)-G-ESTADNQ3-EKREQIYW3-V-N-ELIVM3-E.

NAME: Dihydropteroate synthase signature 1.

- CDZ-x-d-x-ENVIJ-x-EAGJ-ELIVHJ-(2)-N-x-T-x-D-Z-F-x-D-x-ESGJ-

10 NAME: Dihydropteroate synthase signature 2.

CONSENSUS: EGED-ESAD-x-ELIVMD(2)-D-ELIVMD-G-EGPD-x(2)-ESTAD-

x-P.

NAME: EPSP synthase signature 1.

15 CONSENSUS: ELIVMU-x(2)-ENDU-N-ENDU-T-D-EAU-T-ESTAU-x-R-x-ELIVMYU-x-

NAME: EPSP synthase signature 2.

CONSENSUS: EKRJ-x-EKHJ-E-ECZJJ-EDNEJ-R-ELIVMJ-x-ESTAJ-

20 ELIVMCJ-x(2)-EENJ-ELIVMFJ-x-CONSENSUS: EKRAJ-ELIVMFJ-G.

NAME: FLAP/GST2/LTC4S family signature.

G-x(3)-F-B-V-EYJ-x-A-ENQJ-x-N-C

25

NAME: Aminotransferases class-I pyridoxal-phosphate attachment site.

CONSENSUS: EGSI-ELIVMFYTACI-EGSTAI-K-x(2)-EGSALVNI-ELIVMFAI-

x-EGNAR3-x-R-ELIVMA3-

30 CONSENSUS: EGAI.

NAME: Aminotransferases class-II pyridoxal-phosphate

attachment site.

CONSENSUS: T-ELIVMFYWJ-ESAGJ-K-ESAGJ-ELIVMFYWRJ-ESAGJ-x(2)-

35 ESAGI.

NAME: Aminotransferases class-III pyridoxal-phosphate attachment site.

CONSENSUS: ELIVMFYWCJ(2)~x-D-E-ELIVMAJ-x(2)-EGPJ-x(0,1)-

40 ELIVMFYWAGI-x(D,1)-ESACRI-x-

- (E)x-X-EAZDJ-(E/S)x-ECUVYAMVIJJ-C-(d/L-SL)x-ECAZJ-(2,3)-EVZDJ-EVZDJ-(3)-

NAME: Aminotransferases class-IV signature.

45 CONSENSUS: E-x-ESTAGCIJ-x(2)-N-ELIVMFACJ-EFYJ-x(6-12)-

CLIVMF3-x-T-x(6-8)-CLIVM3-x-

CONSENSUS: EGSJ-ELIVMJ-x-EKRJ.

NAME: Aminotransferases class-V pyridoxal-phosphate

50 attachment site.

CONZENZUZ: ELIVFYCHTJ-EHDGJ-ELIVMFYACJ-ELIVMFYAJ-x(2)-

EGSTACI-EGSTAI-EHQRI-K-

CONSENSUS: x(4,6)-G-x-EGSATJ-x-ELIVMFYSACJ.

55 NAME: Hexokinases signature.

CONSENSUS: ELIVMD-G-F-ETND-F-S-EFYD-P-x(5)-ELIVMD-EDNSTD-

 \times (3)-ELIVMD- \times (2)-W-T-K- \times -

CONSENSUS: ELFI.



NAME: Galactokinase signature.

CONSENSUS: G-R-x-N-ELIVI-I-G-E-H-x-D-Y.

5 NAME: GHMP kinases putative ATP-binding domain.

CONSENSUS: ELIVAD-EPKD-x-EGTAD-x(O,1)-G-L-EGSD-2-Z-EGSAD-

-ESATZDB

NAME: Phosphofructokinase signature.

10 CONSENSUS: ERKJ-x(4)-G-H-x-Q-EQRJ-G-G-x(5)-D-R.

NAME: pfkB family of carbohydrate kinases signature 1.
CONSENSUS: EAGI-G-x(O₁1)-EGAPI-x-N-x-ESTAI-x(b)-EG2I-x(9)-G.

15 NAME: pfkB family of carbohydrate kinases signature 2.

CONSENSUS: CDNSKI-EPSTVI-x-CSAGI(2)-EGDI-D-x(3)-ESAGVI-CAGIELIVMYII-ELIVMSTAPI.

NAME: ROK family signature.
20 CONSENSUS: LIVMI-x(2)-G-ELIVMFCTI-G-x-EGAI-ELIVMFAI-x(8)-G-

x(3,5)-EQTADJ-x(2)-CONSCUS: G-EKHJ-

NAME: Phosphoribulokinase signature.

25 CONSENSUS: K-ELIVMI-x-R-D-x(3)-R-G-x-ESTI-x-E.

NAME: Thymidine kinase cellular-type signature.

CONSENSUS: EGAD-x-L3-x-L3-x-L3-x-L3-x-ECHD-

ELIVMFYWHD.

NAME: FGGY family of carbohydrate kinases signature lCONSENSUS: EMFYGSI-x-EPSTI-x(2)-K-ELIVMFYWI-x-W-ELIVMFI-x-

EDENGTKRI-EENGHI.

35 NAME: FGGY family of carbohydrate kinases signature 2-CONSENSUS: EGSAD-x-ELIVMFYWD-x-G-ELIVMD-x(7,8)-EHDENQD-ELIVMFD-x(2)-EASD-ESTAIVMD-

CONSENSUS: ELIVMFY3-EDEQ3.

CONSENSUS: x(5,18)-ELIVMFYWCSTARI-EAIVPI-ELIVMFAGCKRI-K.

45 NAME: Serine/Threonine protein kinases active-site signature.

CONSENSUS: CLIVMFYCD-x-CHYD-x-D-CLIVMFYD-K-x(2)-N-CLIVMFYCTD(3).

50 NAME: Tyrosine protein kinases specific active-site signature.
CONSENSUS: ELIVMFYCI-x-EHYI-x-D-ELIVMFYI-ERSTACI-x(2)-N-

ELIVMFYC3(3).

55 NAME: Protein kinase domain profile.

NAME: Casein kinase II regulatory subunit signature.

CONSENSUS: C-P-x-ELIVMYD-x-C-x(5)-L-P-ELIVMCB-G-x(9)-V-EKRD-x(2)-C-P-x-C

NAME: Pyruvate kinase active site signature.

5 CONSENSUS: ELIVAC3-x-ELIVM3(2)-ESAPCV3-K-ELIV3-E-ENKRST3-x-EDEQH3-EGSTA3-ELIVM3.

NAME: Shikimate kinase signature.

CONSENSUS: EKRD-x(2)-E-x(3)-ELIVMFD-x(8-12)-ELIVMFD(2)-ESAD-

10 x-G(3)-x-ELIVMF3.

40

55

NAME: Prokaryotic diacylglycerol kinase signature.
CONSENSUS: E-x-ELIVMJ-N-ESTJ-ESAJ-ELIVJ-E-x(2)-V-D.

15 NAME: Phosphatidylinositol 3- and 4-kinases signature 1.

CONSENSUS: ELIVMFACD-K-x(1,3)-EDEAD-EDED-ELIVMCD-R-Q-EDEDx(4)-Q.

NAME: Phosphatidylinositol 3- and 4-kinases signature 220 CONSENSUS: EGSI-x-EAVI-x(3)-ELIVMI-x(2)-EFYHI-ELIVMI(2)-xELIVMFI-x-D-R-H-x(2)-N.

NAME: Acetate and butyrate kinases family signature l-CONSENSUS: ELIVMI(2)-x-ELIVMI-N-x-G-S-ESTI-S-x-EKEI.

NAME: Acetate and butyrate kinases family signature 2.

CONSENSUS: ELIVMAI(2)-x(2)-H-x-G-x-ESTI-ELIVMI-x-EAVIx(3)-G.

30 NAME: Phosphoglycerate kinase signature.
CONSENSUS: EKRHGTCVI-EVTI-ELIVMFI-ELIVMCI-R-x-D-x-N-ESACVI-P.

35
NAME: Glutamate 5-kinase signature.
CONSENSUS: ECSTN3-x(2)-G-x-G-ECJ-EIMJ-x-ESTAJ-K-ELIVMJ-x-ESAJ-ETCAJ-x(2)-EGLVJCONSENSUS: x(3)-G.

NAME: ATP:guanido phosphotransferases active site.
CONSENSUS: C-P-x(O₁1)-ESTI-N-EILI-G-T.

NAME: PTS HPR component histidine phosphorylation site signature.
CONSENSUS: G-ELIVMI-H-ESTAI-R-EPAI-EGSTAI-ESTAMI.

NAME: PTS HPR component serine phosphorylation site signature.

50 CONSENSUS: EGSADEJ-EKREQTVJ-x(4)-EKRNJ-S-ELIVMFJ(2)-x-ELIVMJx(2)-ELIVMJ-EGADJ.

NAME: PTS EIIA domains phosphorylation site signature l. CONSENSUS: G-x(2)-ELIVMFl(3)-H-ELIVMFl-G-ELIVMFl-x-T-EALVl.

NAME: PTS EIIA domains phosphorylation site signature 2.
CONSENSUS: EDENQD-x(b)-ELIVMFD-EGAD-x(2)-ELIVMD-A-ELIVMD-P-HEGACD.

NAME: PTS EIIB domains cysteine phosphorylation site

signature.

CONSENSUS: N-ELIVMFJ-x(5)-C-x-T-R-ELIVMFJ-x-ELIVMFJ-x-

5 ELIVMD-x-EDQD.

NAME: Adenylate kinase signature.

CONSENSUS: CLIVMFYWD(3)-D-G-EFYID-P-R-x(3)-ENQD.

10 NAME: Nucleoside diphosphate kinases active site.

CONSENSUS: N-x(2)-H-EGAJ-S-D-ESAJ-ELIVMPKNEJ.

NAME: Guanylate kinase signature.

CONSENSUS: T-EST3-R-x(2)-EKR3-x(2)-EDE3-x(2)-G-x(2)-Y-x-EFY3-

15 ELIVMKI.

NAME: Guanylate kinase domain profile.

NAME: Phosphoribosyl pyrophosphate synthetase signature.

20 CONSENSUS: D-ELID-H-ESAD-x-Q-EIMSTD-EQMD-G-EFYD-F-x(2)-P-

ELIVMFCI-D.

NAME: 7-8-dihydro-b-hydroxymethylpterin-pyrophosphokinase

signature.

25 CONSENSUS: G-EPED-R-x(2)-D-L-D-ELIVMD(2).

NAME: Bacteriophage-type RNA polymerase family active site

signature 1.

CONSENSUS: P-ELIVM3-x(2)-D-EGA3-EST3-EAC3-EGA3-ELIVMFY3-

30. Q.

NAME: Bacteriophage-type RNA polymerase family active site

signature 2.

CONSENSUS: CLIVMFJ-x-R-x(3)-K-x(2)-CLIVMFJ-M-CPTJ-x(2)-Y.

.35

NAME: Eukaryotic RNA polymerase II heptapeptide repeat.

. CANATZI---Z-ETZI--Y-ETZI-Y

NAME: RNA polymerases beta chain signature.

40 CONSENSUS: G-x-K-ELIVMFAD-ESTACD-EGSTND-x-EHSTAD-EGSD-EQNHD-

K-G-EIVT3.

NAME: RNA polymerases M / 15 Kd subunits signature.

CONSENSUS: F-C-x-EDEKSTJ-C-EGNKJ-EDNSAJ-ELIVMHJ-ELIVMJ-

45 x(8,14)-C-x(2)-C.

NAME: RNA polymerases D / 3D to 4D Kd subunits signature.

CITAD-x-3-ENGD-E-x-ELID-

50 CONSENSUS: EGAD-x-R-ELID-EGAD-ELIVMD(2)-P.

NAME: RNA polymerases H / 23 Kd subunits signature.

CONSENSUS: H-ENEID-ELIVMD-V-P-x-H-x(2)-ELIVMD-x(2)-EDED.

55 NAME: RNA polymerases K / 14 to 18 Kd subunits signature.

Q.

NAME: RNA polymerases L / l3 to l6 Kd subunits signature.

CONSENSUS: EDED(2)-H-ESTD-ELIVMD-EGAPD-N-x(l1)-V-x-EFMD-x(2)-Y-x(3)-H-P.

5 NAME: RNA polymerases N / A Kd subunits signature-CONSENSUS: ELIVMFI(2)-P-ELIVMI-x-C-F-ESTI-C-G-

NAME: DNA polymerase family A signature.

CONSENSUS: R-x(2)-EGADAH-X(3)-ELIGNOSHEYD-EAGAD-x(2)-Y-x(2)-

10 EGSD-x(3)-ELIVMAD.

NAME: DNA polymerase family B signature.

CONSENSUS: EYAD-EGLIVMSTACD-D-T-D-ESGD-ELIVMFTCD-xELIVMSTACD.

15

NAME: DNA polymerase family X signature.

CONSENSUS: G-ESGI-ELFYI-x-R-EGEI-x(3)-ESGCLI-x-D-ELIVMI-D-ELIVMI-D-ELIVMI-X(2)-ESAPI.

20 NAME: Galactose-l-phosphate uridyl transferase family lactive site signature.

CONSENSUS: F-E-N-ERK1-G-x(3)-G-x(4)-H-P-H-x-Q.

NAME: Galactose-1-phosphate uridyl transferase family 2

25 signature. CONSENSUS: D-L-P-I-V-G-G-ESTI-ELIVMI(2)-ESAI-H-EDENI-H-EFYI-Q-G-G.

NAME: ADP-glucose pyrophosphorylase signature 1.

30 CONSENSUS: EAG1-G-G-x-G-ESTK1-x-L-x(2)-L-ETA1-x(3)-A-x-P-A-ELV1.

NAME: ADP-glucose pyrophosphorylase signature 2.
CONSENSUS: W-EFY3-x-G-EST3-A-EDNSH3-EAS3-ELIVMFYW3.

35

50

40 NAME: Phosphatidate cytidylyltransferase signature.

CONSENSUS: S-x-ELIVMFI-K-R-x(4)-K-D-x-EGSAI-x(2)-ELII-EPGI-x-H-G-G-ELIVMI-x-D-R
CONSENSUS: ELIVMFII-D.

45 NAME: Ribonuclease PH signature.

CONSENSUS: C-EDEI-ELIVMI(2)-Q-EGTAI-D-G-ESGI-x(2)-ETAI-A.

NAME: 2'-5'-oligoadenylate synthetases signature l.
CONSENSUS: G-G-S-x-EAGI-EKRI-x-T-x-L-EKRI-EGSTI-x-S-D-EAGI-

NAME: 2'-5'-oligoadenylate synthetases signature 2.
CONSENSUS: R-P-V-I-L-D-P-x-EDEI-P-T.

NAME: CDP-alcohol phosphatidyltransferases signature-55 CONSENSUS: D-G-x(2)-A-R-x(3)-G-x(3)-D-x(3)-D.

NAME: PEP-utilizing enzymes phosphorylation site signature.

CONSENSUS: G-EGAD-x-ENTD-x-H-ESTAD-ESTAVD-ELIVMD(2)-ESTAVD-

ERGI.

NAME: PEP-utilizing enzymes signature 2.

5 CONSENSUS: EDERSI-x-ELIVMFI-S-ELIVMFI-G-ESTI-N-D-ELIVMI-x-Q-

LLIVMFYGTJ-ESTALIVJ-

CONSENSUS: CLIVMF3-CGAS3-x(2)-R.

NAME: Rhodanese signature 1.

10 CONSENSUS: EFYD-x(3)-H-ELIVD-P-G-A-x(2)-ELIVFD.

NAME: Rhodanese C-terminal signature.

CONSENSUS: EAVD-x(2)-EFYD-EDEAPD-G-EGSAD-EUFD-x-E-EFYWD.

15 NAME: CoA transferases signature 1.

-q-x-D-(E)x-q-D-D-D-(E) CATMVIJJ-(C)x-CNDJ-CNCJ

NAME: CoA transferases signature 2.

CONSENSUS: CLFJ-CHQJ-S-E-N-G-CLIVFJ(2)-CGAJ.

20

NAME: Phospholipase A2 histidine active site.

CONZENZUZ: C-C-x(5)-H-x(5)-C---

NAME: Phospholipase A2 aspartic acid active site.

25 CONSENSUS: ELIVMAD-C-{LIVMFYWPCST}-C-D-x(5)-C-

NAME: Lipases, serine active site.

CONSENSUS: ELIVI-x-ELIVFYI-ELIVMSTI-G-EHYWVI-S-x-G-EGSTACI.

30 NAME: Colipase signature.

CONSENSUS: Y-x(2)-Y-Y-x-C-x-C

NAME: Lipolytic enzymes "G-D-S-L" family, serine active

site.

35 CONSENSUS: ELIVMFYAGI(4)-G-D-S-ELIVMI-x(1,2)-ETAGI-G.

NAME: Lipolytic enzymes "G-D-X-G" family, putative histidine

active site.

-ENGTZJ-(E)-x-ELYNJ-B-B-B-H-G-B-ESAGJ-EFYJ-x(3)-EZNGJ-ENGTZJ-

40 x(2)-EST3-H.

NAME: Lipolytic enzymes "G-D-X-G" family, putative serine

active site.

CONSENSUS: ELIVMI-x-ELIVMFI-ESAI-G-D-S-ECAI-G-EGAI-x-L-ECAI.

AS

NAME: Carboxylesterases type-B serine active site.

CONSENSUS: F-EGRI-G-x(4)-ELIVMD-x-ELIVD-x-G-x-S-ESTAGD-G.

NAME: Carboxylesterases type-B signature 2.

50 CONSENSUS: CEDI-D-C-L-CYTI-CLIVI-CDNSI-CLIVI-CLIVFYWI-x-

EPQRI.

NAME: Pectinesterase signature 1.

CONSENSUS: EGSTNI-x(5)-ELIVMI-x-ELIVMI-x(2)-G-x-Y-EDNKI-E-x-

55 ELIVMD-x-ELIVMD.

NAME: Pectinesterase signature 2.

CONSENSUS: G-ESTADJ-ELIVMTJ-D-F-I-F-G.

NAME: Peptidyl-tRNA hydrolase signature 1.

CONSENSUS: EFYD-x(2)-T-R-H-N-x-G-x(2)-ELIVMFAD(2)-EDED.

5 NAME: Peptidyl-tRNA hydrolase signature 2.

CONSENSUS: EGSJ-x(3)-H-N-G-ELIVMJ-EKRJ-EDNSJ-ELIVMTJ.

NAME: Alkaline phosphatase active site.

.T-EADJ-ETZADJ-EZZADJ-EZADJ-Z-G-x-EVIJ :ZUZNJZNOJ

NAME: Histidine acid phosphatases phosphohistidine

signature.

-x-R-x-ENDJ-H-R-x-EMVIJJ-(2)x-EMVIJJ-(2)x-ENJJ-x-R-x-

EPASJ.

NAME: Histidine acid phosphatases active site signature.
CONSENSUS: CLIVMFI-x-CLIVMFAGI-x(2)-ESTAGII-H-D-ESTANQI-x-

LIVMJ-x(2)-ELIVMFYJ-x(2)CATZJ

20

NAME: Class A bacterial acid phosphatases signature. CONSENSUS: G-S-Y-P-S-G-H-T.

NAME: 5'-nucleotidase signature 1.

25 CONSENSUS: ELIVMJ-x-ELIVMJ(2)-EHEAJ-ETIJ-x-D-x-H-EGSAJ-x-

CLIVMFD.

NAME: 5'-nucleotidase signature 2.

30

NAME: Fructose-1-6-bisphosphatase active site.

CONSENSUS: EAGI-ERKI-L-x(1,2)-ELIVI-EFYI-E-x(2)-P-ELIVMI-

EGSAJ.

35 NAME: Serine/threonine specific protein phosphatases

signature.

CONSENSUS: ELIVMI-R-G-N-H-E.

NAME: Protein phosphatase 2A regulatory subunit PR55

40 signature 1.

CONZENSUS: E-F-D-Y-L-K-S-L-E-I-E-E-K-I-N.

NAME: Protein phosphatase 2A regulatory subunit PR55

signature 2.

45 CONSENSUS: N-EAGI-H-ETAI-Y-H-I-N-S-I-S-ELIVMI-N-S-D.

NAME: Protein phosphatase 2C signature.

CONSENSUS: ELIVATODELIVATODE CONSENSUS ELIVATODE ELIVATODE CONSENSUS ELIVATODE ELIVATO

EGAVI.

NAME: Tyrosine specific protein phosphatases active site.

CONSENSUS: ELIVMFI-H-C-x(2)-G-x(3)-ESTCI-ESTAGPI-x-ELIVMFYI.

NAME: Tyrosine specific protein phosphatases profile.

55

50

NAME: Dual specificity protein phosphatase profile.

NAME: PTP type protein phosphatase profile.

NAME: Inositol monophosphatase family signature L.

CONSENSUS: EFWVD-x(D-L)-ELIVMD-D-P-ELIVMD-D-ESGD-ESTD-x(2)
EFYD-x-EHKRNSTYD.

5
NAME: Inositol monophosphatase family signature 2.
CONSENSUS: EWVJ-D-x-EACJ-EGSAJ-EGSAVJ-x-ELIVACPJ-ELIVJ-.
ELIVACJ-x(3)-EGHJ-EGAJ.

10 NAME: Prokaryotic zinc-dependent phospholipase C signature.
CONSENSUS: H-Y-x-EGTI-D-ELIVMI-EDNSI-x-P-x-H-EPAI-x-N.

NAME: Phosphatidylinositol-specific phospholipase X-box domain profile.

NAME: Phosphatidylinositol-specific phospholipase Y-box domain profile.

NAME: 3'5'-cyclic nucleotide phosphodiesterases signature20 CONSENSUS: H-D-ELIVMFYD-x-H-x-EAGB-x(2)-ENQD-x-ELIVMFYD-

NAME: cAMP phosphodiesterases class-II signature.

CONSENSUS: H-x-H-L-D-H-ELIVMI-x-EGSI-ELIVMAI-ELIVMI(2)-x-SEAPI.

NAME: Sulfatases signature 1.

CONSENSUS: ESAPI-ELIVMSTI-ECSI-ESTAG-P-ESTAG-R-x(2)ELIVMFW1(2)-ETR1-G.

30 NAME: Sulfatases signature 2.
CONSENSUS: G-EYVI-x-ESTI-x(2)-EIVAI-G-K-x(0,1)-EFYWKI-EHLI.

NAME: AP endonucleases family 1 signature 1.
CONSENSUS: EAPFI-D-ELIVMF1(2)-x-ELIVM1-Q-E-x-K.

35
NAME: AP endonucleases family 1 signature 2.
CONSENSUS: D-ESTI-EFYN-R-EKHI-x(7-8)-EFYNJ-ESTI-EFYWJ(2).

NAME: AP endonucleases family 1 signature 3.
40 CONSENSUS: N-x-G-x-R-ELIVMI-D-ELIVMFYHI-x-ELVI-x-S.

NAME: AP endonucleases family 2 signature 1. CONSENSUS: H-x(2)-Y-ELIVMFJ-EIMJ-N-ELIVMCAJ-EAGJ.

45 NAME: AP endonucleases family 2 signature 2. CONSENSUS: EGRI-ELIVMFI-C-ELIVMI-D-T-C-H.

NAME: Deoxyribonuclease I signature l.

CONSENSUS: ELIVMI(2)-EAPI-L-H-ESTAI(2)-P-x(5)-E-ELIVMI-EDNIx-L-x-EDEI-V.

55 NAME: Deoxyribonuclease I signature 2.
CONSENSUS: G-D-F-N-A-x-C-ESAI.

NAME: Endonuclease III iron-sulfur binding region signature.

CONSENSUS: C-x(3)-EKRSI-P-EKRAGLI-C-x(2)-C-x(5)-C.

NAME: Endonuclease III family signature.

-CZAGD-CIJD-(C, S)x-CWMVIJD-(2)x-G-C7MVID-x-CTZDD : ZUZN3ZNO)

5 G-V-EGAD-x(3)-EGACJ-

CONZENZUZ: x(3)-ELIVIJ-x(2)-EZALVJ-ELIVMYJ-EGANKJ.

NAME: Ribonuclease II family signature.

CONSENSUS: EHID-EFYED-EGSTAMD-ELIVMD-x(4,5)-Y-ESTALD-x-

10 EFWVACD-ETVD-ESAD-P-ELIVMAD-

CONSENSUS: CRQ3-CKR3-CFY3-x-D-x(3)-CHQ3.

NAME: Ribonuclease III family signature.

CONSENSUS: EDEQI-ERQI-ELMI-E-EFYWI-ELVI-G-D-ESARI.

15

NAME: Bacterial Ribonuclease P protein component signature. CONSENSUS: ELIVMFYSJ-x(2)-A-x(2)-R-ENHJ-EKRQLJ-ELIVMJ-EKRAJ-R-x-ELIVMTAJ-EKRJ.

V-X-FFTAILIVR-FKVT.

20 NAME: Ribonuclease T2 family histidine active site 1.

NAME: Ribonuclease T2 family histidine active site 2.

CONSENSUS: CLIVMF3-x(2)-EHDGTYJ-EEQ3-EFYWJ-x-EKRJ-H-G-x-C.

25

NAME: Pancreatic ribonuclease family signature.

CONSENSUS: C-K-x(2)-N-T-F.

NAME: DNA/RNA non-specific endonucleases active site.

30 CONSENSUS: D-R-G-H-EQILI-x(3)-A.

NAME: Thermonuclease family signature 1.

CONSENSUS: D-G-D-T-CLIVMD-x-CLIVMCD-x(9,10)-R-CLIVMD-x(2)-

ELIVMI-D-x-P-E.

35

NAME: Thermonuclease family signature 2.

CONSENSUS: D-EKRJ-Y-EGQJ-R-x-ELVJ-EGAJ-x-EIVJ-EFYWJ.

NAME: Beta-amylase active site 1.

40 CONSENSUS: H-x-C-G-G-N-V-G-D.

NAME: Beta-amylase active site 2.

CONSENSUS: G-x-ESAB-G-E-ELIVMB-R-Y-P-S-Y.

45 NAME: Glucoamylase active site region signature.

CONSENSUS: ESTNI-EGPI-x(L,2)-EDEI-x-U-E-E-x(2)-EGSI.

NAME: Polygalacturonase active site.

CONSENSUS: EGSDENKHID-x(2)-EVHFCD-x(2)-EGSD-H-G-ELIVMAGD-

 $50 \times (1-2)-ELIVMI-G-S$.

NAME: Clostridium cellulosome enzymes repeated domain

signature.

CONSENSUS: D-ELIVMFJJ-EVNGJ-x-EVNGJ-x-(2)-ELIVMJ-B-CZALMJ-

55 x-D-x(3)-ELIVMF3-x-

NAME: Chitinases family 18 active site.

CUNZENZUZ: ELIVATIJENDI-G-ELIVAFIJ-EDNI-ELIVAFIJ-EDNI-x-E.

NAME: Chitinases family 19 signature 1.

CONSENSUS: $(E)_{x-A-x} = (E)_{x-A-x} = (E)$

 $5 \times (2) - F - \mathbb{E}GSAI$

NAME: Chitinases family 19 signature 2.

CONSENSUS: ELIVMJ-ECSAJ-F-x-ESTAGJ(2)-ELIVMFYJ-W-EFYJ-W-

ELIVMI.

10

NAME: Alpha-lactalbumin / lysozyme C signature.

C- \times (3)- \mathbb{C} - \times (3)- \mathbb{C} - \times (3)- \mathbb{C} DENJ- \mathbb{C} LIJ- \times (5)- \mathbb{C} .

NAME: Alpha-galactosidase signature.

15 CONSENSUS: G-ELIVMFYD-x(2)-ELIVMFYD-x-ELIVMD-D-x-W-x(3,4)-

R-EDNSFI.

NAME: Trehalase signature 1.

CONSENSUS: P-G-G-R-F-x-E-x-Y-x-W-D-x-Y.

20

NAME: Trehalase signature 2.

CONSENSUS: Q-W-D-x-P-x-EGAD-W-EPAD-P.

NAME: Alpha-L-fucosidase putative active site.

25 CONSENSUS: P-x(2)-L-x(3)-K-W-E-x-C.

NAME: Glycosyl hydrolases family 1 active site.

CONSENSUS: ELİYMYZJ-ELYYVJ-ELLYBYZJ-ELYZHWYJJ-E-N-G-

ELIVMFARD-ECSAGND.

30

NAME: Glycosyl hydrolases family 1 N-terminal signature.
CONSENSUS: F-x-EFYWMJ-EGSTAJ-x-EGSTAJ-x-EGSTAJ(2)-EFYNHJ-

ENQ3-x-E-x-EGSTA3.

35 NAME: Glycosyl hydrolases family 2 signature 1.

CONSENSUS: N-x-ELIVMFYWDD-R-ESTACND(2)-H-Y-P-x(4)-

ELIVATION (E) x (2) = CDND
40 NAME: Glycosyl hydrolases family 2 acid/base catalyst.

CONSENSUS: EDENGED-EKRVWD-H-K-EDD-ESCD-ELIVMFD(3)-W-EGSD-

.3-N-E.

NAME: Glycosyl hydrolases family 3 active site.

45 CONSENSUS: ELIVMI(2)-EKRI-x-EEQKI-x(4)-G-ELIVMFTI-ELIVTI-

ELIVMFI-ESTI-D-x(2)-

CONSENSUS: ESGADNII.

NAME: Glycosyl hydrolases family 5 signature.

50 CONSENSUS: CLIVI-ELIVMFYWGAI(2)-EDNEQGI-ELIVMGSTI-x-N-E-EPVI-

ERHDNSTLIVEY3.

NAME: Glycosyl hydrolases family b signature 1.

CONSENSUS: V-x-Y-x(2)-P-x-R-D-C-EGSAFJ-x(2)-EGSAJ(2)-x-G.

55

NAME: Glycosyl hydrolases family & signature 2.

CONSENSUS: CLIVMYAD-CLIVAD-CLIVAD-CLIVA-E-P-D-CSALD-CLID-

EPSAGI.

NAME: Glycosyl hydrolases family & signature.

CONSENSUS: A-ESTI-D-EAGI-D-x(2)-EIMI-A-x-ESAI-ELIVMI-ELIVMGI-x-A-x(3)-EFWI.

NAME: Glycosyl hydrolases family 9 active sites signature 1-CONSENSUS: ESTV3-x-ELIVMFY3-ESTV3-x(2)-G-x-ENKR3-x(4)-EPLIVM3-H-x-R.

10 NAME: Glycosyl hydrolases family 9 active sites signature 2-CONSENSUS: EFYWD-x-D-x(4)-EWYDD-x(3)-E-x-EATAD-x(3)-N-EZTAD.

NAME: Glycosyl hydrolases family 10 active site.

CONSENSUS: EGTA1-x(2)-ELIVN1-x-EIVMF1-EST1-E-ELIY1-EDN1-

15 ELIVMFI.

NAME: Glycosyl hydrolases family ll active site signature l. CONSENSUS: EPSAU-ELQU-x-E-Y-Y-ELIVMU(2)-EDEU-x-EFYWHNU.

20 NAME: Glycosyl hydrolases family ll active site signature 2. CONSENSUS: ELIVMFD-x(2)-E-EAGD-EYWGD-EQRFGSD-ESGD-ESTAND-G-x-ESAFD.

NAME: Glycosyl hydrolases family 16 active sites.

25 CONSENSUS: E-ELIVI-D-ELIVI-x(0,1)-E-x(2)-EGQI-EKRNFI-x-EPSTAI.

NAME: Glycosyl hydrolases family 17 signature.

CONSENSUS: ELIVMI-x-ELIVMVII)-ESTAGI-E-ESTAI-G-W-P-ESTNI-

30 x-ESAGQI.

NAME: Glycosyl hydrolases family 25 active sites signature.
CONSENSUS: D-ELIVMI-x(3)-ENQI-EPGI-x(9,10)-G-x(4)ELIVMFYI(2)-K-x-ESTI-E-EGSI-x(2)-

35 CONSENSUS: Y-x-EDNJ.

NAME: Glycosyl hydrolases family 31 active site. CONSENSUS: EGFI-ELIVMFI-W-x-D-M-ENSAI-E.

40 NAME: Glycosyl hydrolases family 3L signature 2.

CONSENSUS: G-EAVN-D-ELIVMTN-C-G-EFYN-x(3)-ESTN-x(3)-L-C-x-R-U-x(2)-ELVN-EGSN-ESANCONSENSUS: F-x-P-F-x-R-EDNN.

45 NAME: Glycosyl hydrolases family 32 active site. CONSENSUS: H-x(2)-P-x(4)-ELIVMI-N-D-P-N-G.

NAME: Glycosyl hydrolases family 35 putative active site.
CONSENSUS: G-G-P-ELIVM3(2)-x(2)-Q-x-E-N-E-EFY3.

NAME: Glycosyl hydrolases family 39 active site.
CONSENSUS: W-x-F-E-x-W-N-E-P-ENDD.

NAME: Glycosyl hydrolases family 45 active site.
55 CONSENSUS: ESTAD-T-R-Y-EFYWD-D-x(5)-ECAD.

NAME: Prokarvotic transglycosylases signature.

CMVZIJ= $D-Q-(5)\times -Z-X(3)$

CANSENSUS.

CONSENSUS: x(4)-ESAGI.

5 NAME: Inosine-uridine preferring nucleoside hydrolase family

signature.

CONSENSUS: D-x-D-ETTJ-EGAJ-x-D-D-ETAVJ-EVIJ-A.

NAME: Alkylbase DNA glycosidases alkA family signature.

10 CONSENSUS: G-I-G-x-W-ESTD-EAVD-x-ELIVMFYD(2)-x-ELIVMD-x(8)-

EMF3-x(2)-EED3-D.

NAME: Formamidopyrimidine-DNA glycosylase signature.

-x-EATZJ-ENATZDJ-R-E(7)x-EVIJ-x-CVIZ-x-C(4c)

NAME: Uracil-DNA glycosylase signature.

CONSENSUS: EKRI-ELIVI-ELIVCI-ELIVMI-x-G-EQII-D-P-Y.

20 NAME: S-adenosyl-L-homocysteine hydrolase signature 1.

CONSENSUS: CCSJ-N-x-EFYLJ-S-ESTJ-EQAJ-EDENJ-x-EAVJ(2)-A-A-

-EVAZI-EVIJI

NAME: S-adenosyl-L-homocysteine hydrolase signature 2.

25 CONSENSUS: $G-K-x(3)-\mathbb{E}LIVI-x-G-Y-G-x-V-G-\mathbb{E}KRI-G-x-A$

NAME: Cytosol aminopeptidase signature.

CONSENSUS: N-T-D-A-E-G-R-L.

30 NAME: Aminopeptidase P and proline dipeptidase signature.

CONSENSUS: EHAD-EGYZD-ELIVTD-EGD-H-x-ELIVD-G-ELIVTD-x-

CIVD-H-EDED.

NAME: Methionine aminopeptidase subfamily 1 signature-

35 CONSENSUS: EMFYD-x-G-H-G-ELIVMOD-EGSHD-x(3)-H-x(4)-ELIVMD-x-

EHND-EYWVD.

NAME: Methionine aminopeptidase subfamily 2 signature-

CONSENSU: EDAB-ELIVMYB-x-K-ELIVMB-D-x-G-x-EH@B-ELIVMB-EDNSD-

40 G-x(3)-EDN3-

NAME: Renal dipeptidase active site.

CONSENSUS: ELIVMD-E-G-EGAD-x(2)-ELIVMFD-x(b)-L-x(3)-Y-x(2)-G-

ELIVMI-R.

45

NAME: Serine carboxypeptidases, serine active site.

CONSENSUS: ELIVMD-x-EGTAD-E-S-Y-EAGD-EGSD.

NAME: Serine carboxypeptidases, histidine active site.

50 CONSENSUS: ELIVFID-x(2)-ELIVSTAI-x-EIVPSTI-x-EGSDNQLI-ESAGVI-

ESGU-H-x-EIQAVIU-x-H-EDZU-

CONSENSUS: EPSAJ.

NAME: Zinc carboxypeptidases, zinc-binding region 1

55 signature.

CONSENSUS: EPKI-x-ELIVMFYI-x-ELIVMFYI-x(4)-H-ESTAGI-x-E-x-

ELIVMI-ESTAGI-x(b)-

CONSENSUS: ELIVMFYTAD.

NAME: Zinc carboxypeptidases, zinc-binding region 2 signature.

CONSENSUS: H-ESTAGI-x(3)-ELIVMEI-x(2)-ELIVMFYWI-P-EFYWI.

NAME: Serine proteases, trypsin family, histidine active site.

CONSENSUS: ELIVMI-ESTI-A-ESTAGI-H-C.

10 NAME: Serine proteases, trypsin family, serine active siteCONSENSUS: CONSENSUS: CHAPAMICATEGES - CONSENSUS: CHAPAMICATEGES - CONSENSUS: CHAPAMICATEGES - CONSENSUS: CLIVATEGES - CONSENSUS: CONSENSUS - CONSE

15 NAME: Serine proteases, subtilase family, aspartic acid active site.

CONSENSUS: ESTAIVI-x-ELIVMFI-ELIVMI-D-EDSTAI-G-ELIVMFCIx(2,3)-EDNHI.

20 NAME: Serine proteases, subtilase family, histidine active site.

CONSENSUS: H-G-ESTMI-x-EVICI-ESTAGCI-EGSI-x-ELIVMAIESTAGCLVI-ESAGMI.

25 NAME: Serine proteases, subtilase family, serine active site.

CONSENSUS: G-T-S-x-ESAI-x-P-x(2)-ESTAVCI-EAGI.

NAME: Serine proteases, V8 family, histidine active site.
30 CONSENSUS: ESTI-G-ELIVMFYWI(3)-EGNI-x(2)-T-ELIVMI-x-T-x(2)-H.

NAME: Serine proteases, VB family, serine active site.

CONSENSUS: T-x(2)-EGCB-ENQB-S-G-S-x-ELIVMB-EFYB.

35 NAME: Serine proteases, omptin family signature L. CONSENSUS: W-T-D-x-S-x-H-P-x-T.

40

45

NAME: Serine proteases, omptin family signature 2.
CONSENSUS: A-G-Y-Q-E-ESTI-R-EFYWI-S-EFYWI-ETNI-A-x-G-G-ESTIY.

NAME: Prolyl endopeptidase family serine active site-CONSENSUS: D-x(3)-A-x(3)-ELIVMFYUJ-x(L4)-G-x-S-x-G-G-ELIVMFYUJ(2).

NAME: Endopeptidase Clp serine active site.
CONSENSUS: T-x(2)-ELIVMFJ-G-x-A-ESACJ-S-EMSAJ-EPAGJ-ESTAJ.

NAME: Endopeptidase Clp histidine active site.
50 CONSENSUS: R-x(3)-EEAJ-x(3)-ELIVMTYI-M-ELIVMJ-H-Q-P.

NAME: ATP-dependent serine proteases, lon family, serine active site.

CONSENSUS: D-G-EPDJ-S-A-EGSJ-ELIVMCAJ-ELIVMJ.

NAME: Eukaryotic thiol (cysteine) proteases cysteine active site.

CONSENSUS: Q-x(3)-EGEI-x-C-EYWI-x(2)-ESTAGCI-ESTAGCVI.

Eukaryotic thiol (cysteine) proteases histidine active NAME:

site.

CONZENZUZ:

ELIVMGSTAND-x-H-EGSACED-ELIVMD-x-ELIVMATD(2)-G-x-

EGSADNHI.

Eukaryotic thiol (cysteine) proteases asparagine

active site. CONZENZUZ:

EFYCHD-EWID-ELIVTD-x-EKRQAGD-N-ESTD-W-x(3)-EFYWB-

G-x(2)-G-ELFYWII-

CONSENSUS: ELIVMFYGD-x-ELIVMFD.

Ubiquitin carboxyl-terminal hydrolase family 1

cysteine active-site.

CONSENSUS: Q-x(3)-N-ESA3-C-g-x(3)-ELIVN3(2)-H-ESA3-ELIVN3-15

-EAZB

Ubiquitin carboxyl-terminal hydrolases family 2 NAME:

signature 1.

G-ELIVMFYD-x(1,3)-EAGCD-ENASMD-x-C-EFYWD-ELIVMCD-CONSENSUS:

-CZMVIJ-x-CLIVMZJ-

CONSENSUS:

NAME: Ubiquitin carboxyl-terminal hydrolases family 2

25 signature 2.

Y-x-L-x-ESAGI-ELIVMFTI-x(2)-H-x-G-x(4,5)-G-H-Y-CONSENSUS:

Caspase family histidine active site.

H-x(2,4)-ESCU-x(4)-ELIVMFU(2)-ESTU-H-G. CONZENZUZ:

30

20

Caspase family cysteine active site.

CONZENZUZ: K-P-K-ELIVMFJ(4)-Q-A-C-ERQGJ-G.

Eukaryotic and viral aspartyl proteases active site.

35 -- CONSENSUS: ELIVMFGACD-ELIVMTADND-ELIVFSAD-D-ESTD-G-ESTAVD-

ESTAPDENGD-x-ELIVMFSTNCD-

x-ELIVMFGTAI. CONSENSUS:

NAME: Neutral zinc metallopeptidases, zinc-binding region

40 signature.

> [GSTALIVN]-x(2)-H-E-[LIVMFYW]-{DEHRKP}-H-x-CONZENZUZ:

ELIVMFYWGSPQJ.

Matrixins cysteine switch.

CONSENSUS: P-R-C-EGNJ-x-P-EDRJ-ELIVSAPKQJ. 45

Insulinase family, zinc-binding region signature. NAME:

G-x(B-9)-G-x-ESTAD-H-ELIVMFYD-ELIVMCD-EDERND-CONSENSUS:

CHRKLJ-CLMFATJ-x-CLFSTHJ-x-

EGSTANI-EGSTI. 50 CONZENZUZ:

11

FACOTOSH AC

55 Glycoprotease family signature.

EKRI-EGSATI-x(4)-EFYWHLI-EDQNGKI-x-P-x-ELIVMFYI-CONZENZUZ:

x(3)-H-x(2)-EAGJ-H-

CONZENZUZ: ILIVMD.



NAME: Proteasome A-type subunits signature.

CONSENSUS: $\mathbb{E} FYJ-x(4)-\mathbb{E}STNVJ-x-\mathbb{E}FYUJ-S-P-x-G-\mathbb{E}KHJ-x(2)-Q-$

LLIVMD-EDEJ-Y-EGADJ-x(2)-

5 CONSENSUS: ESAGI.

NAME: Proteasome B-type subunits signature.

CONSENSUS: ELIVMAI-EGZAI-ELIVMFI-x-EFYLVGACI-x(2)-EGZACFYI-

ELIVMSTACI(3)-EGACI-

10 CONSENSUS: EGSTACVJ-EDESJ-x(15)-ERKJ-x(12,13)-G-x(2)-EGSTAJ-D.

NAME: Signal peptidases I serine active site.

CONSENSUS: EGS3-x-S-M-x-EPS3-EAT3-ELF3.

15

NAME: Signal peptidases I lysine active site.

CONSENSUS: K-R-ELIVM3-AJ(2)-6-x-EPGJ-G-EDEJ-x-ELIVM3-x-

ELIVMFYD.

20 NAME: Signal peptidases I signature 3.

CONSENSUS: $\mathbb{C}LIVMFYWJ(2)-x(2)-G-D-\mathbb{C}HJ-x(3)-\mathbb{C}NDJ-x(2)-\mathbb{C}GJ$.

NAME: Signal peptidases II signature.

CYRYVIJ-EGAMVJA-CASJ-EHVIJ-EGASJ-N-ELYMYGJ-ELADJ-ELAZJ-N-ELYMYGJ-ELIVRYJ-

25 D-R-ELIMFAI.

NAME: Peptidase family U32 signature.

CONSENSUS: E-x-F-x(2)-G-ESAJ-ELIVMJ-C-x(4)-G-x-C-x-ELIVMJ-S.

30 NAME: Amidases signature.

-EAZƏJ-EMVIJ-x-EYVAZƏJ-EAZƏJ-x-Ə-EZƏJ-Z-Z-EAƏJ- SUZNJZNOJ

x(b)-EGSAI-x-EGAI-x-D-

CONSENSUS: x-EGAD-x-S-ELIVMD-R-x-P-EGSACD.

35 NAME: Asparaginase / glutaminase active site signature l.

CONSENSUS: ELIVMD-x(2)-T-G-G-T-EIVD-EAGSD.

NAME: Asparaginase / glutaminase active site signature 2.

CONSENSUS: G-x-ELIVMD-x(2)-H-G-T-D-T-ELIVMD.

40

NAME: Urease nickel ligands signature.

NAME: Urease active site.

45 CONSENSUS: ELIVMD(2)-ECTD-H-EHND-L-x(3)-ELIVMD-x(2)-D-ELIVMD-

x-F-A.

NAME: ArgE / dapE / ACYL / CPG2 / yscS family signature l.

CONSENSUS: ELIVI-EGALMYI-ELIVMFI-x-EGSAI-H-x-D-ETVI-ESTAVI.

50

NAME: ArgE / dapE / ACY1 / CPG2 / yscS familý signature 2.

CLIVMFYD-x(14-17)-ELIVMD-

CONSENSUS: X-ELIVMFJ-ELIVMSTAGJ-ELIVMFAJ-x(2)-EDNGJ-E-E-x-

55 EGSTNI.

NAME: Dihydroorotase signature 1.

CONSENSUS: D-ELIVMFYWSAPJ-H-ELIVAJ-H-ELIVFJ-ERNJ-x-EPGNJ.



Dihydroorotase signature 2. NAME: CONSENSUS: EGAD-ESTD-D-x-A-P-H-x(4)-K.

NAME: Beta-lactamase class-A active site.

EFYD-x-ELIVMFYD-x-S-ETVD-x-K-x(4)-EAGLMD-x(2)-CONSENSUS:

ELCI.

NAME: Beta-lactamase class-C active site. CONZENZUZ: F-E-CLIVMI-G-S-CLIVMGI-CSAI-K. 10

NAME: Beta-lactamase class-D active site.

CONSENSUS: EPAD-x-S-ESTD-F-K-ELIVD-EPALD-x-ESTAD-ELID.

15 Beta-lactamases class B signature 1.

ELID-x-ESTND-ENND-x-H-EGSTAD-D-x(2)-G-EGPD-x(7,8)-CONZENZUZ:

EGSJ.

Beta-lactamases class B signature 2.

P-x(3)-ELIVMJ(2)-x-G-x-C-ELIVMFJ(2)-K. 20 CONZENZUZ:

Arginase family signature 1.

CONSENSUS: ELIVMFJ-G-G-x-H-x-ELIVMTJ-ESTAVJ-x-EPAGJ-x(3)-

EGSTAD.

25

35

40

50

Arginase family signature 2. NAME:

CONZENZUZ: CLIVMJ(2)-x-ELIVMFYJ-D-EASJ-H-x-D.

Arginase family signature 3. NAME:

CONSENSUS: -CAZQJ-9-(E)x-CQA9-(E)x-C-CMVIJJ-D-CY7MVIJJ-CTZJ 30

x(7)-G.

NAMF: Adenosine and AMP deaminase signature.

CONZENSUS: EZAJ-ELIVMJ-ENGSJ-ESTAJ-D-D-P-

NAME: Cytidine and deoxycytidylate deaminases zinc-binding

region signature.

CONZENZUZ: ECHI-EAGVI-E-x(2)-ELIVMFGATI-ELIVMI-x(17,33)-P-C-

 $\times (2 - B) - C - \times (3) - ELIVM3$

GTP cyclohydrolase I signature 1. NAME:

CONSENSUS: CEND-CLIVMD(2)-x(2)-CKRQND-CND-CLIVMD-x(3)-CSTD-

x-C-E-H-H-

45 NAME: GTP cyclohydrolase I signature 2.

CONSENSUS: ESAB-x-ERKD-x-Q-ELIVMD-Q-E-ERND-ELID-ETSND.

NAME: Nitrilases / cyanide hydratase signature 1.

CONZENZUZ: G-x(2)-CLIVMFYJ(2)-x-CIFJ-x-E-x(2)-CLIVMJ-x-G-Y-P.

Nitrilases / cyanide hydratase active site signature. NAME: **CONSENSUS:**

EKRI-

55 Inorganic pyrophosphatase signature.

CONZENZUZ: D-ESGDND-D-EPED-ELIVMFD-D-ELIVMGACD.

NAME: Acylphosphatase signature 1. 5

PCT/IB01/02050

CONSENSUS: ELIVI-x-G-x-V-Q-G-V-x-EFMI-R.

NAME: Acylphosphatase signature 2.

CONSENSUS: G-EFYUJ-EAVCJ-EKRQAMJ-N-x(3)-G-x-V-x(5)-G.

NAME: ATP synthase alpha and beta subunits signature.

CONSENSUS: P-EFAZJ-EVIJ-EVIJ-EVIJ-S-x-Z.

NAME: ATP synthase gamma subunit signature.

10 CONSENSUS: EIVI-T-x-E-x(2)-GII-x(3)-G-A-x-ESAKRI.

NAME: ATP synthase delta (OSCP) subunit signature.

CONSENSUS: ELIVMD-x-ELIVMFYTD-x(3)-ELIVMTD-EDENQKD-x(2)-

ELIVMD-x-EGSAD-G-ELIVMFYGAD-

15 CONSENSUS: x-ELIVMI-EKRHENQI-x-EGSENI.

20 NAME: ATP synthase c subunit signature.

CONSENSUS: EGSTAl-R-ENQl-P-x(l0)-ELIVMFYWl(2)-x(3)-ELIVMFYWl-x-EDEl.

NAME: E1-E2 ATPases phosphorylation site.

25 CONSENSUS: D-K-T-G-T-ELID-ETID.

NAME: Sodium and potassium ATPases beta subunits signature

CONSENSUS: EFYWB-x(2)-EFYWB-x-EFYWB-EDND-x(b)-ELIVMB-G-R-T-

30 x(3)-W-

NAME: Sodium and potassium ATPases beta subunits signature

2.

CONSENSUS: ERK3-x(2)-C-ERKQWI3-x(5)-L-x(2)-C-ESA3-G.

NAME: GDAL/CD39 family of nucleoside phosphatases signature.

CONSENSUS: ELIVMJ-x-G-x(2)-E-G-x-EFYJ-x-EFWJ-ELIVAJ-ETAGJ-x-N-EHYJ.

40 NAME: Iodothyronine deiodinases active site-CONSENSUS: R-P-L-V-x-N-F-G-S-ECAI-T-C-P-x-F-

NAME: Cutinase, serine active site.

CONSENSUS: P-x-ESTAD-x-ELVD-EIVTD-x-EGSD-G-Y-S-EQLD-G-

NAME: DDC / GAD / HDC / TyrDC pyridoxal-phosphate attachment

50 site.

CONSENSUS: S-ELIVMFYWJ-x(5)-K-ELIVMFYWGJ(2)-x(3)-ELIVMFYWJ-x-ECAJ-x(2)-ELIVMFYWZJ
CONSENSUS: x(2)-ERKJ.

NAME: Orn/DAP/Arg decarboxylases family 2 pyridoxal-P

attachment site.

CONSENSUS: EFYJ-EPAJ-x-K-ESACVJ-ENHCLFWJ-x(4)-ELIVMFJ-

5 ELIVMIJA (2) - ELIVMIJA (3) - ESTEJ : SUZNAZNO)

NAME: Orn/DAP/Arg decarboxylases family 2 signature 2.
CONSENSUS: EGS3-x(2,b)-ELIVMSCP3-x(2)-ELIVMF3-EDNS3-ELIVMCA3-

10 G-G-G-ELIVMFYI-

CONSENSUS: EGSTPCEQI.

NAME: Orotidine 5'-phosphate decarboxylase active site-CONSENSUS: ELIVMFTAD-ELIVMFD-x-D-x-K-x(2)-D-I-EGPD-x-T-

15 ELIVMTAD.

50

NAME: Phosphoenolpyruvate carboxylase active site 1.
CONSENSUS: EVT3-x-T-A-H-P-T-EEQ3-x(2)-R-EKRH3.

20 NAME: Phosphoenolpyruvate carboxylase active site 2. CONSENSUS: EIV3-M-ELIVM3-G-Y-S-D-S-x-K-D-ESTAG3-G.

NAME: Phosphoenolpyruvate carboxykinase (GTP) signature. CONSENSUS: F-P-S-A-C-G-K-T-N.

NAME: Phosphoenolpyruvate carboxykinase (ATP) signature.
CONSENSUS: L-I-G-D-D-E-H-x-W-x-EDEJ-x-G-EIVJ-x-N.

NAME: Uroporphyrinogen decarboxylase signature 1-30 CONSENSUS: P-x-W-x-M-R-Q-A-G-R.

NAME: Uroporphyrinogen decarboxylase signature 2.

CONSENSUS: G-F-ESTAGCV3-ESTAGC3-x-P-EFYW3-T-ELV3-x(2)-Y-x(2)EAE3-EGK3-

35
NAME: Indole-3-glycerol phosphate synthase signature.
CONSUS: ELIVATY3-ELIVACI-x-E-ELIVATYCI-K-EKRSPI-ESTAKI-S-P-ETZI-x(3)-ELIVATVII.

40 NAME: Ribulose bisphosphate carboxylase large chain active site.
CONSENSUS: G-x-EDNI-F-x-K-x-D-E.

NAME: Fructose-bisphosphate aldolase class-I active site45 CONSENSUS: CLIVMI-x-CLIVMFYWI-E-G-x-CLSI-L-K-P-CSNI.

NAME: Malate synthase signature.

55 CONSENSUS: EKRI-EDENGI-H-x(2)-G-L-N-x-G-x-W-D-Y-ELIVMI-F.

NAME: Hydroxymethylglutaryl-coenzyme A lyase active site-CONSENSUS: S-V-A-G-L-G-G-C-P-Y.

NAME: Hydroxymethylglutaryl-coenzyme A synthase active site-CONSENSUS: N-x-EDNJ-EIVJ-E-G-EIVJ-D-x(2)-N-A-C-EFYJ-x-G.

5 NAME: Citrate synthase signature.
CONSENSUS: G-EFYAI-EGAI-H-x-EIVI-x(1,2)-ERKTI-x(2)-D-EPSI-R.

NAME: Alpha-isopropylmalate and homocitrate synthases signature 1.

10 CONSENSUS: L-R-EDED-G-x-Q-x(LD)-K.

NAME: Alpha-isopropylmalate and homocitrate synthases signature 2.

NAME: KDPG and KHG aldolases active site-CONSENSUS: G-ELIVMD-x(3)-E-ELIVD-T-ELFD-R.

NAME: KDPG and KHG aldolases Schiff-base forming residue. 20 CONSENSUS: G-x(3)-ELIVMTJ-K-ELFJ-F-P-ESAJ-x(3)-G.

NAME: Isocitrate lyase signature.
CONSENSUS: K-EKRI-C-G-H-ELMQI.

25 NAME: Beta-eliminating lyases pyridoxal-phosphate attachment site.

CONSENSUS: Y-x-D-x(3)-M-S-EGAD-K-K-D-x-ELIVMD(2)-x-ELIVMD-G-

30 NAME: DNA photolyases class & signature &.

CONSENSUS: T-G-x-P-ELIVMJ(2)-D-A-x-M-ERAJ-x-ELIVMJ.

NAME: DNA photolyases class 1 signature 2.

CONSENSUS: EDN1-R-x-R-ELIVM1(2)-x-ESTA1(2)-F-ELIVMFA1-x-K-x35 L-x(2,3)-W-EKRQ1.

NAME: DNA photolyases class 2 signature 1.
CONSENSUS: F-x-E-E-x-ELIVM1(2)-R-R-E-L-x(2)-N-F.

40 NAME: DNA photolyases class 2 signature 2.

CONSENSUS: G-x-H-D-x(2)-W-x-E-R-x-ELIVM3-F-G-K-ELIVM3-R-EFY3-M-N.

NAME: Eukaryotic-type carbonic anhydrases signature.
45 CONSENSUS: S-E-H-x-ELIVM3-x(4)-EFYH3-x(2)-E-ELIVM3-H-ELIVMFA3(2).

NAME: Prokaryotic-type carbonic anhydrases signature 1.
CONSENSUS: C-ESAB-D-S-R-ELIVMB-x-EAPB.

55 NAME: Fumarate lyases signature. CONSENSUS: G-S-x(2)-M-x(2)-K-x-N.

NAME: Aconitase family signature 1.

CONSENSUS: ELIVAD.

5 NAME: Aconitase family signature 2.
CONSENSUS: G-x(2)-ELIVWPQI-x(3)-EGACI-C-EGSTAMI-ELIMPTAI-C-

ELIMVI-EGAI.

NAME: Dihydroxy-acid and L-phosphogluconate dehydratases

10 signature 1.

 $CAD = CE \times CAD = CE$

NAME: Dihydroxy-acid and b-phosphogluconate dehydratases signature 2.

15 CONSENSUS: ESAD-L-ELIVMD-T-D-EGAD-R-ELIVMFD-S-EGAD-EGAVD-ESTD.

NAME: Dehydroquinase class I active site.

CONZENZUZ: D-ELIVMJ-EDEJ-ELIVNJ-x(L8-20)-ELIVMJ(2)-x-ESCJ-

20 ENHYD-H-EDND.

EFYD-G.

25

NAME: Enolase signature.

CONSENSUS: ELIVI(3)-K-x-N-Q-I-G-ESTI-ELIVI-ESTI-EDEI-ESTAI.

NAME: Serine/threonine dehydratases pyridoxal-phosphate attachment site.

CONSENSUS: EDESHB-x(4.5)-ESTVGD-x-EASB-EFYID-K-EDLIFSAB-ERVMFD-EGAD-ELIVMGAD.

NAME: Enoyl-CoA hydratase/isomerase signature.

35 CONSENSUS: : CATZĂRANADAU-EMVIJU : ZUZNAZNO) - (E) EDAU-(E) X-Q-(ETZĀRĀNĀNADAU-(H) X-(E) EDAU-X-(E) ZUZNAZNO) : ZUZNAZNO) : ZUZNAZNO)

NAME: Imidazoleglycerol-phosphate dehydratase signature 1.
40 CONSENSUS: ELIVMYD-EDED-x-H-H-x(2)-E-x(2)-EGCAD-ELIVMDESTACD-ELIVMD.

NAME: Imidazoleglycerol-phosphate dehydratase signature 2. CONSENSUS: G-x-EDNI-x-H-H-x(2)-E-ESTAGCI-x-EFYI-K.

45
NAME: Tryptophan synthase alpha chain signature.
CONSENSUS: ELIVMI-E-ELIVMI-G-x(2)-EFYCI-ESTI-EDEI-EPABELIVMYI-EAGLII-EDEI-G.

50 NAME: Tryptophan synthase beta chain pyridoxal-phosphate attachment site.

CONSENSUS: ELIVMI-x-H-x-G-ESTAI-H-K-x-N.

NAME: Delta-aminolevulinic acid dehydratase active site.
55 CONSENSUS: G-x-D-x-ELIVMI(2)-EIVI-K-P-EGSAI-x(2)-Y.

NAME: Urocanase active site.
CONSENSUS: F-Q-G-L-P-x-R-I-C-W.

NAME: Prephenate dehydratase signature 1.

CLIVMWD-x-ELIVMB.

5

NAME: Dihydrodipicolinate synthetase signature 1.

10 CONSENSUS: EGSAJ-ELIVMJ-ELIVMTJ-x(2)-G-ESTJ-ETGJ-G-E-EGJ-x(b)-EEQJ.

NAME: Dihydrodipicolinate synthetase signature 2.

-(41-EL)x-EMVIJJ-(E)x-ET2J-(S)x-Q-E7MVIJJ-EZMCJ-Y

15 CLIVMI-x-CSGAB-CLIVMFI-

CONSENSUS: K-EDEQAFI-ESTACI.

NAME: RsuA family of pseudouridine synthase signature. CONSENSUS: G-R-L-D-x(2)-ESTI-x-G-ELIVMFI(4)-ESTI-EDNTI.

20

NAME: Cysteine synthase/cystathionine beta-synthase Pphosphate-attachment site.

CONSENSUS: K-x-E-x(3)-EPAD-ESTAGCD-x-S-EIVAPD-K-x-R-x-ESTAGDx(2)-ELIVMD.

25

NAME: Phenylalanine and histidine ammonia-lyases signature.
CONSENSUS: G-ESTGJ-ELIVMJ-ESTGJ-ECJ-S-G-EDHJ-L-x-P-L-ESAJx(2)-ESAJ.

30 NAME: Porphobilinogen deaminase cofactor-binding site.
CONSENSUS: E-R-x-ELIVMFAJ-x(3)-ELIVMFJ-x-G-EGSAJ-C-x-EIVTJ-PELIVMFJ-EGSAJ.

NAME: Cys/Met metabolism enzymes pyridoxal-phosphate
35 attachment site.
CONSENSUS: EDQ3-ELIVMF3-x(3)-ESTAGC3-ESTAGC13-T-K-EFYWQ3ELIVMF3-x-G-EHQ3-ESGNH3.

NAME: Glyoxalase I signature 1.

40 CONSENSUS: EHQ3-EIVI3-x-ELIVI3-x-EIVI3-x-(2)-F-EYM3-x(2,3)-ELMF3-G-ELMF3.

NAME: Glyoxalase I signature 2.

CONSENSUS: G-ENTKQD-x(D,5)-EQAD-ELYFYD-EGHD-H-EIVFD-ECAD-x-

45 ESTAGLU-x(2)-EDNCU.

NAME: Cytochrome c and cl heme lyases signature l-CONSENSUS: H-N-x(2)-N-E-x(2)-W-ENQKRI-x(4)-W-E.

50 NAME: Cytochrome c and cl heme lyases signature 2. CONSENSUS: P-F-D-R-H-D-W.

NAME: Adenylate cyclases class-I signature l-CONSENSUS: E-Y-F-G-ESAI(2)-L-W-x-L-Y-K-

55

NAME: Adenylate cyclases class-I signature 2.
CONSENSUS: Y-R-N-x-W-ENSI-E-ELIVMI-R-T-L-H-F-x-G.

NAME: Guanylate cyclases signature.

CONSENSUS: G-V-ELIVM3-x(0,1)-G-x(5)-EFY3-x-ELIVM3-EFYW3-EGS3-

EDNTHKWD-EDNTD-EIVD-

CONSENSUS: CDNTAI-x(5)-CDEI.

5

NAME: Chorismate synthase signature 1.

CONZENSUS: G-E-S-H-LGCJ-x(2)-LLIVMJ-EGTVJ-x-LLIVMJ(2)-LDEJ-G-

x-EPVI.

10 NAME: Chorismate synthase signature 2.

CONSENSUS: EGEJ-R-ESAJ(2)-ESAJJ-R-EVJ-ESTJ-x(2)-EHJ-V-x(2)-

G.

NAME: Chorismate synthase signature 3.

15 CONSENSUS: R-ESHD-D-EPSVD-ECSAVD-x(4)-EGAID-x-EIVGSPD-ELIVMD-

X-E-ESTAHJ-ELIVMJ.

NAME: 6-pyruvoyl tetrahydropterin synthase signature 1.

CONSENSUS: C-N-N-x(2)-G-H-G-H-N-Y.

20

30

NAME: 6-pyruvoyl tetrahydropterin synthase signature 2.

CONSENSUS: D-H-K-N-L-D-x-D.

NAME: Ferrochelatase signature.

25 CONSENSUS: ELIVMFJ(2)-x-S-x-H-EGSJ-ELIVMJ-P-x(4,5)-EDENQKRJ-

x-G-D-x-Y.

NAME: Alanine racemase pyridoxal-phosphate attachment site.

CONSENSUS: V-x-K-A-EDN3-EGA3-Y-G-H-G.

NAME: Aspartate and glutamate racemases signature l-CONSENSUS: CIVAL-CLIVML-x-C-x(0-1)-N-CSTL-CMSAL-CSTHL-

ELIVFYSTANKJ.

35 NAME: Aspartate and glutamate racemases signature 2.

CONSENSUS: ELIVMI(2)-x-EAGI-C-T-EDEHI-ELIVMFYI-EPNGRSI-x-

ELIVMI.

NAME: Mandelate racemase / muconate lactonizing enzyme

40 family signature 1.

CONSENSUS: A-x-ESAGI(2)-ELIVMI-EDEI-x-A-x(2)-D-x(2)-EGAI-

EKR3.

NAME: Mandelate racemase / muconate lactonizing enzyme

45 family signature 2.

CONSENSUS: G-x(7)-D-x(9)-A-x(14)-ELIMID-E-EDENQI-P-x(4)-

EDENGI.

NAME: Ribulose-phosphate 3-epimerase family signature 1.

50 CONSENSUS: ELIVMFJ-H-ELIVMFYJ-D-ELIVMJ-x-D-x(1-2)-EFYJ-

-EVAT23-x-N-x-EMVIJ3

NAME: Ribulose-phosphate 3-epimerase family signature 2.

CONSENSUS: ELIVMAJ-x-ELIVMJ-M-EZTJ-EVZJ-x-P-x(3)-G-Q-x-F-

55 x(b)-ENK3-ELIVMC3-

NAME: Aldose 1-epimerase putative active site.

CONSENSUS: ENSI-x-T-N-H-x-Y-EFWI-N-ELII.

NAME: Cyclophilin-type peptidyl-prolyl cis-trans isomerase signature.

CONSENSUS: EFYD-x(2)-ESTCNLVD-x-F-H-ERHD-ELIVMD-x(2)-F-ELIVMD-x-Q-EAGD-G.

NAME: Cyclophilin-type peptidyl-prolyl cis-trans isomerase profile.

- 10 NAME: FKBP-type peptidyl-prolyl cis-trans isomerase signature 1.

 CONSENSUS: ELIVMCI-x-EYFI-x-EGVLI-x(1,2)-ELFTI-x(2)-G-x(3)-EDEI-ESTAERKI-ESTANI.
- NAME: FKBP-type peptidyl-prolyl cis-trans isomerase domain profile.
- NAME: PpiC-type peptidyl-prolyl cis-trans isomerase

 25 signature.

 CONSENSUS: F-EGSADEII-x-ELVAQI-A-x(3)-ESTI-x(3,4)-ESTQIx(3,5)-EGERI-G-x-ELIVMICONSENSUS: EGSI-
- 30 NAME: Triosephosphate isomerase active site. CONSENSUS: EAVI-Y-E-P-ELIVMI-W-ESAI-I-G-T-EGKI.
 - NAME: Xylose isomerase signature L. CONSENSUS: ELII-E-P-K-P-x(2)-P.
- NAME: Xylose isomerase signature 2.
 CONSENSUS: EFLI-H-D-x-D-ELIVI-x-EPDI-x-EGDEI.
- NAME: Phosphomannose isomerase type I signature 1-40 CONSENSUS: Y-x-D-x-N-H-K-P-E-
 - NAME: Phosphomannose isomerase type I signature 2.

 CONSENSUS: H-A-Y-ELIVMI-x-G-x(2)-ELIVMI-E-x-M-A-x-S-D-N-x-ELIVMI-R-A-G-x-T-P-K.
- 45
 NAME: Phosphoglucose isomerase signature l.
 CONSENSUS: EDENSU-x-ELIVMU-G-G-R-EFYU-S-ELIVMU-x-ESTAUEPSACU-ELIVMU-G.
- 50 NAME: Phosphoglucose isomerase signature 2.

 CONSENSUS: EGSI-x-ELIVMI-ELIVMFYWI-x(4)-EFYI-EDNI-Q-x-G-V-Ex(2)-K.
- NAME: Glucosamine/galactosamine-b-phosphate isomerases
 55 signature.
 CONSENSUS: FLIVID-x(3)-G-x-FLITID-x-FLIVID-x-FLIVID-x-G-FLIVID-

NAME: Phosphoglycerate mutase family phosphohistidine

signature.

CONSENSUS: ELIVMI-x-R-H-G-EEQI-x(3)-N.

5 NAME: Phosphoglucomutase and phosphomannomutase

phosphoserine signature.

CONSENSUS: EASAJ-ELTVIJ-x-ELTVIJ-ERSAJ-S-H-x-P-x(4)-

EGNHEJ-

10 NAME: Methylmalonyl-CoA mutase signature.

CONSENSUS: R-I-A-R-N-ETQ3-x(2)-ELIVMFY3(2)-x-EEQ3-E-x(4)-

EKRN3-x(2)-D-P-x-EGSA3-

CONSENSUS: G-S.

15 NAME: Terpene synthases signature.

CONSENSUS: EDED-G-S-W-x-G-x-W-EGAD-ELIVMD-x-EFYD-x-Y-EGAD.

NAME: Eukaryotic DNA topoisomerase I active site.

CENTRAL STREET S

20 ELIVME.

NAME: Prokaryotic DNA topoisomerase I active site.

EDEQSI.

25

NAME: DNA topoisomerase II signature-

CONSENSUS: CLIVMAD-x-E-G-EDND-S-A-x-ESTAGD-

NAME: Aminoacyl-transfer RNA synthetases class-I signature.

30 CONSENSUS: P-x(D-2)-EGSTAND-EDENGGAPKD-x-ELIVMFPD-EHTD-

ELIVMYACD-G-EHNTGD-

CONSENSUS: ELIVMFYSTAGPCI-

NAME: Aminoacyl-transfer RNA synthetases class-II signature

35 1.-

CONSENSUS: CFYHJ-x-x-CDCJ-x(4,12)-EHJ-x(3)-F-x(3)-EDEJ.

NAME: Aminoacyl-transfer RNA synthetases class-II signature

2.

40 CONSENSUS: EGSTALVFI-{DANGHRKP}-EGSTAI-ELIVHFI-EDEI-R-

ELIVMF3-x-ELIVMSTAG3-ELIVMFY3.

NAME: WHEP-TRS domain signature.

CONSENSUS: EQUI-G-ENGUI-X-ELIVI-EKRI-X(2)-K-X(2)-EKRNGI-

45 EASI-x(4)-ELIVI-EDENKI-

CONZENZUZ: $\times(5)-\mathbb{L}\Lambda\mathbb{J}-\times(5)-\Gamma-\times(3)-K$

NAME: ATP-citrate lyase / succinyl-CoA ligases family

signature 1.

50 CONSENSUS: S-EKRI-S-G-EGTI-ELIVMI-EGSTI-x-EEQI-x(8,10)-G-

x(4)-ELIVMI-EGAI-ELIVMI-G-

CONSENSUS: G-D.

NAME: ATP-citrate lyase / succinvl-CoA ligases family active

55 site.

CONSENSUS: $G-x(2)-A-x(4-7)-\mathbb{E}RQTJ-\mathbb{E}LIVMFJ-G-H-\mathbb{E}ASJ-\mathbb{E}GHJ$.

NAME: ATP-citrate lyase / succinyl-CoA ligases family

signature 3.

-CVATZJ-Q-CLADJ-Q-CLANJ-x-CRMVIJ-x-CVIJ-x-CVLD-x-G-CLAJ-CZVA-Z

 $-(E)\times-\mathbb{E}MVIJJ-\times-(E)\times$

5 CONSENSUS: G-EGREI-

NAME: Glutamine synthetase signature 1.

CDNSCHOLD = CENTRAL SALES OF THE CONSCHOOL OF THE CONSCHO

[LIVMFY].

10

NAME: Glutamine synthetase putative ATP-binding region

signature.

-(E)x-H-x-D-ENAT2DJ-D-ETA9NJ-(3,5)-ENPATJ-G-EGZNAJ-G-X-H-x(3)-

- 7

15

NAME: Glutamine synthetase class-I adenylation site.

CONSENSUS: K-ELIVMJ-x(5)-ELIVMAJ-D-ERKJ-EDNJ-ELIJ-Y.

NAME: D-alanine--D-alanine ligase signature l.

20 CONSENSUS: H-G-x(2)-G-E-D-G-x-ELIVMAD-EQSAD-EGSAD.

NAME: D-alamine--D-alamine ligase signature 2.

CONSENSUS: ELIVI-x(3)-ECAD-x-EGSALVI-R-ELIVCAD-D-ELIVMFI(2)-

x(7,9)-ELI3-x-E-

25 CONSENSUS: CLIVAI-N-ESTPI-x-P-EGAI.

NAME: SAICAR synthetase signature 1.

CONSENSUS: ELIVMTJ(2)-P-ELIVMTJ-E-x-ELIVMTJ-ELIVMCAT-R-x(3)-

ETAI-6-5.

30

NAME: SAICAR synthetase signature 2.

CONSENSUS: ELIVMD-ELIVMAD-D-x-K-ELIVMFYD-E-F-G.

NAME: Folylpolyglutamate synthase signature 1.

35 CONSENSUS: ELIVMFYD-x-ELIVMD-ESTAGD-G-T-ENKD-G-K-x-ESTD-x(7)-

ELIVMI(2)-x(3)-EGZKI.

NAME: Folylpolyglutamate synthase signature 2.

CONSENSUS: ELIVMFYD(2)-E-x-G-ELIVMD-EGAD-G-x(2)-D-x-EGSTD-x-

40 [LIVM](2).

NAME: Ubiquitin-activating enzyme signature 1.

CONSENSUS: K-A-C-S-G-K-F-x-P.

45 NAME: Ubiquitin-activating enzyme active site.

CONSENSUS: P-ELIVMI-C-T-ELIVMI-EKRHI-x-EFTI-P.

NAME: Ubiquitin-conjugating enzymes active site.

50 ELIVI-x-ELIVI.

NAME: Formate--tetrahydrofolate ligase signature 1.

CONSENSUS: G-ELIVMI-K-G-G-A-A-G-G-G-Y.

55 NAME: Formate--tetrahydrofolate ligase signature 2.

CONSENSUS: V-A-T-EIVJ-R-A-L-K-x-EHNJ-G-G.

NAME: Adenylosuccinate synthetase GTP-binding site.

CONSENSUS: Q-W-G-D-E-G-K-G.

NAME: Adenylosuccinate synthetase active site.

CONZENZUZ: $G-I-\mathbb{E}GR\mathbb{I}-P-x-Y-x(2)-K-x(2)-R$

5 NAME: Argininosuccinate synthase signature 1.

CONZENZUZ: A-EFYD-S-G-L-D-T-S.

Argininosuccinate synthase signature 2. CONZENZUZ: 10 G-x-T-x-K-G-N-D-x(2)-R-F.

Phosphoribosylglycinamide synthetase signature. NAME:

R-F-G-D-P-E-x-EQMD. CONZENSUS:

15 Carbamoyl-phosphate synthase subdomain signature 1. CONSENSUS: EFYVD-EPSD-ELIVMCD-ELIVMAD-ELIVMD-EKRD-EPSAD--EDAD-x-D-EDZD-(E)x-EATZD

Carbamoyl-phosphate synthase subdomain signature 2. 20 CONZENZUZ: ELIVMFD-ELIMND-E-ELIVMCAD-N-EPATLIVMD-EKRD-

ELIVMSTAC3.

NAME: ATP-dependent DNA ligase AMP-binding site. CONZENZUZ: CEDQH3-x-K-x-CDNJ-G-x-R-CGACIVMJ.

25

ATP-dependent DNA ligase signature 2. NAME:

CONZENSUS: E-G-ELIVMAD-ELIVMD(2)-EKRD-x(5-8)-EYWD-EQNEKD-

x(2-L)-EKRHI-x(3-5)-K-

CONSENSUS: ELIVMFYI-K.

30

NAD-dependent DNA ligase signature 1.

CONSENSUS: K-ELIVMD-D-G-ELIVMD-ESAD-x(4)-Y-x(2)-G-x-L-x(4)-

 $\mathbb{L}ST\mathbb{J}-R-G-\mathbb{L}DN\mathbb{J}-G-x(2)-G-$

CONSENSUS: EDEB-EDENLB.

35

NAD-dependent DNA ligase signature 2. NAME:

EIVI-G-EKRI-ESTI-G-x-ELIVMI-ESTNKI-x-EVTI-x(2)-L-**CONZENZUZ:**

-V-EZGJ-x

40 NAME: RNA 3'-terminal phosphate cyclase signature.

CONZENZUS: $\mathbb{C} \mathbb{R} \mathbb{H} \mathbb{J} - \mathbb{G} - \mathbb{X} (2) - \mathbb{P} - \mathbb{X} - \mathbb{G} (3) - \mathbb{X} - \mathbb{C} \mathbb{L} \mathbb{I} \mathbb{V} \mathbb{J}$

NAME: Lipoate-protein ligase B signature.

R-G-G-x(2)-T-EFYWJ-H-x(2)-EGHJ-Q-x-ELIVJ-x-Y. **CONSENSUS:**

45

NAME: Isopenicillin N synthetase signature 1.

CONZENZUZ: $\mathbb{C}\mathbb{R}\mathbb{K}\mathbb{J}-x-\mathbb{C}\mathbb{S}\mathbb{T}\mathbb{A}\mathbb{J}-x$ (2)-S-x-C-Y- $\mathbb{C}\mathbb{S}\mathbb{L}\mathbb{J}$.

NAME: Isopenicillin N synthetase signature 2.

50 - FDNG1-x-T-(2)x-EDAT21-(2)x-EAT21-D-0-x-(2)= EMVIJ1 CONZENZUZ:

NAME: Site-specific recombinases active site.

CONZENZUS: Y-LLIVACJ-R-LVAJ-S-LSTJ-x(2)-Q.

55 Site-specific recombinases signature 2.

CONZENZUS: G-EDED-x(2)-ELIVMD-x(3)-ELIVMD-EDTD-R-ELIVMD-

EGSAI.

NAME: Transposases, Mutator family, signature.

CONZENZUZ: D-x(3)-G-ELIVMFI-x(b)-ESTAVI-ELIVMFYWI-EPTI-x-

 $-(5)x-0-x-\mathbb{E}RDD-(5)x-\mathbb{E}VAT2D$

CONSENSUS: н.

5 Transposases, ISBO family, signature.

CONSENSUS: R-G-x(2)-E-N-x-N-G-ELIVMI(2)-R-EQEI-ELIVMFYI(2)-P-

10 NAME: Autoinducers synthetases family signature.

CONZENZUZ: CLMFYD-R-x(3)-F-x(2)-CKRD-x(2)-W-x-CLIVMD-x(6,9)-

E-x-D-x-EFYI-D.

Thiamine pyrophosphate enzymes signature. NAME:

CONSENSUS: 15 ELIVMFI-EGSAI-x(5)-P-x(4)-ELIVMFYWI-x-ELIVMFI-x-G-

D-EGSAI-EGSACI.

Biotin-requiring enzymes attachment site.

CONSENSUS: EGND-EDEQTRD-x-ELIVMFYD-x(2)-ELIVMD-x-EAIVD-M-K-

20 ELMATI-x(3)-ELIVMI-x-

CONSENSUS: - EVAZI

2-oxo acid dehydrogenases acyltransferase component

lipoyl binding site.

25 CONZENZUZ: EGND-x(2)-ELIVFD-x(5)-ELIVFCD-x(2)-ELIVFAD-x(3)-K-

-ENGQVATZJ-EVIATZJ

x(2)-ELIVMFSD-x(5)-EGCND-x-ELIVMFYD. CONSENSUS:

Putative AMP-binding domain signature.

30 **CONSENSUS:** -x-EDZJ-ETZJ-ETZJ-ETZJ-EZZJ-EZZJ-EZZJ-EZZJ-ZZ

EPASLIVMI-EKRI.

Molybdenum cofactor biosynthesis proteins signature 1.

CONSENSUS:

ELIVM3(3)-ELIT3(2)-G-G-T-G-x(4)-D.

Molybdenum cofactor biosynthesis proteins signature 2. NAME: . **CONZENZUZ:** $S-x-\mathbb{C}SJ-x(2)-D-x(5)-\mathbb{C}LIVUJ-x(JO,J2)-\mathbb{C}LIVJ-x(2)-$

EKRI-P-G-EKRLI-P-x(2)-

CONZENZUZ: ELIVMFD-EGAD.

40

35

NAME: moaA / nifB / pqqE family signature.

ELIVI-x(3)-C-ENPI-ELIVMFI-EQRSI-C-x-EFYMI-C. CONSENSUS:

Radical activating enzymes signature. NAME:

CONSENSUS: EGV3-x-G-x-EKR3-x(3)-F-x(2)-G-x(0,1)-C-x(3)-C-45

x(2)-C-x-ENLII.

NAME: Tpx family signature.

S-x-D-L-P-F-A-x(2)-EKRJ-EFWJ-C. **CONSENSUS:**

50

NAME: Cytochrome c family heme-binding site signature.

CONZENZUZ: C-{CPWHF}-{CPWR}-C-H-{CFYW}.

Cytochrome b5 family, heme-binding domain signature. NAME:

55 CONSENSUS: EFY3-ELIVMKJ-x(2)-H-P-EGAJ-G.

NAME: Cytochrome b/bb heme-ligand signature.

CONZENSUS: EDENGI-x(3)-G-EFYWMQI-x-ELIVMFI-R-x(2)-H.

NAME: Cytochrome b/bb @o site signature.
CONSENSUS: P-EDEJ-W-EFYJ-ELFYJ(2).

5 NAME: Cytochrome b559 subunits heme-binding site signature. CONSENSUS: ELIVI-x-ESTI-ELIVFI-R-EFYWI-x(2)-EIVI-H-ESTGAI-ELIVI-ESTGAI-EIVI-P.

NAME: Nickel-dependent hydrogenases b-type cytochrome subunit signature 1.

CONSENSUS: R-ELIVMFYWW-x-H-W-ELIVMW-x(2)-ELIVMFW-ESTACW-ELIVMW-x(2)-L-x-ELIVMW-T-G-

NAME: Nickel-dependent hydrogenases b-type cytochrome subunit signature 2.

CONSENSUS: ERH3-ESTAD-ELIVMFYWD-H-ERHD-ELIVMD-x(2)-W-x-ELIVMFD-x(2)-F-x(3)-H.

NAME: Succinate dehydrogenase cytochrome b subunit signature

20 l.
CONSENSUS: R-P-ELIVMTD-x(3)-ELIVMD-x(b)-ELIVMWPKD-x(4)-S-

CONSENSUS: R-P-ELIVMID-x(3)-ELIVMID-x(6)-ELIVMWPKID-x(4)-S-x(2)-H-R-x-ESTI.

NAME: Succinate dehydrogenase cytochrome b subunit signature 25 2.
CONSENSUS: H-x(3)-EGAI-ELIVMTI-R-EHFI-ELIVMFI-x-EFYWMI-D-x-EGVAI.

NAME: Thioredoxin family active site-

30 CONSENSUS: ELIVMFID-ELIVMŠTAID-x-ELIVMFYCID-EFYWSTHEID-x(2)EFYWGTNID-C-EGATPLVEIDCONSENSUS: EPHYWSTAID-C-x(b)-ELIVMFYWII.

NAME: Glutaredoxin active site.

40

35 CONSENSUS: CLIVDI-CFYSAI-x(4)-C-CPVI-CFYUI-C-(2)-CTAVI- x(2-3)-CLIVI-

NAME: Type-l copper (blue) proteins signature.

CONSENSUS: EGAl-x(0,2)-EYSAl-x(0,1)-EVFYl-x-C-x(l,2)-EPGl-x(0,1)-H-x(2,4)-EMQl.

NAME: 2Fe-2S ferredoxins, iron-sulfur binding region signature.

CONSENSUS: C-{C}-{C}-EGAI-{C}-C-EGASTI-{CPDEKRHFYW}-C-

NAME: Adrenodoxin family, iron-sulfur binding region signature.

CONSENSUS: C-x(2)-ESTAQI-x-ESTAMVI-C-ESTAI-T-C-EHRI.

50 NAME: 4Fe-4S ferredoxins, iron-sulfur binding region signature.
CONSENSUS: C-x(2)-C-x(2)-C-x(3)-C-EPEGI.

NAME: High potential iron-sulfur proteins signature.

55 CONSENSUS: C-x(b-9)-ELIVMI-x(3)-G-EYWI-C-x(2)-EFYWI.

NAME: Rieske iron-sulfur protein signature l. CONSENSUS: C-ETKI-H-L-G-C-ELIVII.

NAME: Rieske iron-sulfur protein signature 2. CONSENSUS: C-P-C-H-x-EGSAJ.

Flavodoxin signature.

CONZENZUZ: CLIVD-CLIVFYD-CFYD-x-CSTD-x(2)-CAGCD-x-T-x(3)-Ax(2)-ELIVI-

Rubredoxin signature.

10 CONZENZUZ: ELIVMD-x(3)-W-x-C-P-x-C-EAGDD.

Electron transfer flavoprotein alpha-subunit NAME:

signature.

CONZENZUZ: CLID-Y-CLIVMD-CATD-x-G-CIVD-CSDD-G-x-CIVD-Q-H-

 \times (2)-G- \times (b)- \mathbb{I} IV \mathbb{J} - \times -A-CONZENSUS: LIVI-N-

> Electron transfer flavoprotein beta-subunit signature. NAME:

CONZENZUZ: LIVAB-x-LKRD-x(2)-LDED-LGDD-LGDED-x(1,2)-LEQD-x-

20 ELIV3-x(4)-P-x-ELIVM3(2)-

CONZENSUS: ETACE-

Vertebrate metallothioneins signature.

C-x-C-EQAT2DI-x(2)-C-x-C-x(2)-C-x-K-CONZENZUZ:

25 Ferritin iron-binding regions signature 1.

CONSENSUS: E-x-EKR3-E-x(2)-E-EKR3-ELF3-ELIVMA3-x(2)-Q-N-x-Rx-G-R.

30 NAME: Ferritin iron-binding regions signature 2. **CONZENSUS:** -EYJJ-(2)x-CMJJ-T-LLIJ-T-LLIJ-EZJ-LTMVIJJ-(2)x-C L-x(b)-ELIVMD-EKND.

Bacterioferritin signature.

<M-x-G-x(3)-V-ELIVD-x(2)-ELMD-x(3)-L-x(3)-L.</pre> **CONSENSUS:** 35

Transferrins signature 1.

CONSENSUS: Y-x(O-1)-EVASI-V-EIVACI-EIVAI-EIVAI-ERKSI-

EGDENSA3.

Transferrins signature 2.

CONSENSUS: Y-x-G-A-EFLJ-EKRHNQJ-C-L-x(3,4)-G-EDENQJ-V-EGAJ-

EFYW3.

40

45 NAME: Transferrins signature 3.

CONZENZUZ: EDENGU-EYFU-x-ELYU-L-C-x-EDNU-x(5-8)-ELIVU-x(4-5)-

C-x(2)-A-x(4)-EHQR3-x-

CONZENZUZ: ELIVMFYWD-ELIVMD.

50 NAME: Globins profile.

> Protozoan/cyanobacterial globins signature. NAME:

F-ELFJ-x(5)-G-EPAJ-x(4)-G-EKRAJ-x-ELIVMJ-x(3)-H. **CONZENZUZ:**

55 NAME: Plant hemoglobins signature. CONZENZUZ: $\mathbb{E}SNJ-P-x-L-x(2)-H-A-x(3)-F.$

NAME: Hemerythrins signature.

CONSENSUS: W-L-x-ENQD-H-I-x(3)-D-F-

Arthropod hemocyanins / insect LSPs signature 1. **CONZENSUS:** Y-EFYUJ-x-E-D-ELIVMJ-x(2)-N-x(b)-H-x(3)-P.

5 NAME: Arthropod hemocyanins / insect LSPs signature 2. CONZENZUZ: T-x(2)-R-D-P-x-EFYJ-EFYWJ.

Heavy-metal-associated domain.

CONZENZUZ: 10 ELIVNI-x(2)-ELIVMFAI-x-C-x-ESTAGCDHI-(-x(3)-ELIVFGD-x(3)-ELIVD-x(9,11)-CONZENSUS: CIVAB-x-ELVFYSD.

NAME: ABC transporters family signature.

CONSENSUS: ELIVMFYCD-ESAD-ESAPGLVFYKQHD-G-EDENQMUD-15 EKRQASPCLIMFW3-EKRNQSTAVM3-

CONZENSUS: EKRACLVMJ-ELIVMFYPANJ-{PHY}-ELIVMFWJ-ESAGCLIVPJ-{FYWHP}-{KRHP}-

CONSENSUS: ELIVMFYWSTAD.

20

NAME: Binding-protein-dependent transport systems inner membrane comp. sign. CONZENSUS: ELIVMFYD-x(8)-EEQRD-ESTAGVD-ESTAGD-x(3)-G-

ELIVMFYSTAC3-x(5)-ELIVMFYSTA3-

25 CONSENSUS: ×(4)-ELIVMFY3-EPKR3.

ABC-2 type transport system integral membrane proteins NAME: signature.

CONZENZUZ: -EVIAZDD-x-EJTZDD-EADMIJD-(2)x-EUMIJD-(5)x-ETZMIJD

30 x(b)-ELIMGAI-EPGZNQI-

x(9,12)-P-ELIMFTD-x-EHRSYD-x(5)-ERQD. CONZENZUZ:

Bacterial extracellular solute-binding proteins, family 1 signature.

35 CONZENZUZ: -Y-CYAMVIJ-ELYMVAD-ELYMVAD-ELYMVYD(2)-Y-ENDD-x(3)-ELIVMFD-x-CONSENSUS: **EKNDED.**

Bacterial extracellular solute-binding proteins,

,40 family 3 signature. CONSENSUS: G-EFYILD-EDED-ELIVMTD-EDED-ELIVMFD-x(3)-ELIVMAD-EVAGCI-x(2)-ELIVMAGNI.

Bacterial extracellular solute-binding proteins,

45 family 5 signature.

CONSENSUS: TAGD-x(b,7)-TDNEGD-x(2)-TSTAVED-TLIVMFYWAD-x-

ELIVMFYD-x-ELIVMD-EKRD-

CONSENSUS: EKRHDED-EGDND-ELIVMAD-EKNGSPJ-EFWD.

50 Serum albumin family signature. **CONSENSUS:** EFYJ-x(b)-C-C-x(?)-C-ELFYJ-x(b)-ELIVMFYWJ.

NAME: Transthyretin signature 1.

CONZENZUZ: S-K-C-P-L-M-V-K-V-L-D-EASJ-V-R-G.

55

Transthyretin signature 2.

CONSENSUS: S-P-EFYJ-S-EFYJ-S-T-T-A-ELIVMJ-V-ESTJ-x-P.

NAME: Avidin / Streptavidin family signature.

CONSENSUS: EDENJ-x(2)-ERTAJ-x(2)-V-G-x-EDNJ-x-EFUJ-T-

EKRJ-

5 NAME: Eukaryotic cobalamin-binding proteins signature. CONSENSUS: ESNI-V-D-T-EGAI-A-ELIVMI-A-x-L-A-ELIVMFI-T-C.

NAME: Lipocalin signature.

CONSENSUS: EDENGI-x-EDENGGSTARKI-x(0,2)-EDENGARKI-ELIVFYI-

10 (CP)-G-{C}-W-EYWLRH3-x-CONSENSUS: ELIVMTA3.

NAME: Cytosolic fatty-acid binding proteins signature.

CONSENSUS: EGSAIVKD-x-EFYWD-x-ELIVMFB-x(4)-ENHGD-EFYD-EDED-x-

15 CLIVMTYD-CLIVMD-x(2)-

CONSENSUS: ELIVMAKRI.

NAME: Acyl-CoA-binding protein signature.

-CAT203-Y-CYMVIJ3-(2)x-CTMVIJ3-x-CN3Q3-x-CAT23-9

20 x-EFYJ-K-Q-ESTAJ(2)-x-G.

NAME: LBP / BPI / CETP family signature.

CONSENSUS: CPAD-EGAD-ELIVMCD-x(2)-R-EIVD-ECAD-L-x(3)-L-x(5)-

TEQD-x(4)-TLIVMD-TEQKD-

25 CONSENSUS: x(B)-P.

NAME: Phosphatidylethanolamine-binding protein family

signature.

CONSENSUS: EFYD-x-ELIVMFD(3)-x-EDCD-P-D-x-P-ESND-x(10)-H.

30

NAME: Plant lipid transfer proteins signature.

-x-ELIVMI-EYAMI-EX-CONSULT -x-ELIVMI-x-

ELIVMI-ETZI-x(3)-

CONSENSUS: EDN3-C-x(2)-ELIVM3.

35

NAME: Uteroglobin family signature 1.

CONZENZUS: CEAU-x-C-P-x-ELWVIJ-x(3)-ELIVMJ-EDED-x-

ILIVMFI(2).

40 NAME: Uteroglobin family signature 2.

-CONSENSUS: EDSGI-C(4)-ESNI-x(5)-CDSGI-x(-1-x(2)-SUSNI

C-

NAME: Mitochondrial energy transfer proteins signature.

45 CONSENSUS: P-x-EDEJ-x-ELIVATJ-ERKJ-x-ELRHJ-ELIVMFYJ-EQMAIGVJ.

NAME: Sugar transport proteins signature 1.

CONSENSUS: ELIVATION = ELIVATION = EDATEMVIJ = EDATEMVIJ = EUZASANO = -x-

ELIVMFYWAD-G-R-ERKD-x(4,6)-

50 CONSENSUS: EGSTAI.

NAME: Sugar transport proteins signature 2.

CONSENSUS: ELIVMFID-x-G-ELIVMFAD-x(2)-G-x(A)-ELIFYD-x(2)-EEQD-

x(b)-ERKI.

NAME: LacY family proton/sugar symporters signature 1.

CONSENSUS: $G = \mathbb{E}[IVM](2) - x - D = \mathbb{E}[K] - L - G - L - \mathbb{E}[K](2) - x - \mathbb{E}[IVM](2) - \omega$.

NAME: LacY family proton/sugar symporters signature 2.

CONSENSUS: P-x-ELIVMFJ(2)-N-R-ELIVMJ-G-x-K-N-ESTAJ-ELIVMJ(3).

NAME: PTR2 family proton/oligopeptide symporters signature

L.

COUSTAINTING TO THE SECTION
10 NAME: PTR2 family proton/oligopeptide symporters signature 2.

CONSENSUS: [FYT]=x(2)-[LMFY]-[FYV]-[LIVMFYWA]-x-[IVG]-N[LIVMAG]-G-[GSA]-[LIMF].

15 NAME: Amiloride-sensitive sodium channels signature.

CONSENSUS: Y-x(2)-EEQTFJ-x-C-x(2)-EGZTDNLJ-C-x-EQTJ-x(2)
ELIVITJ-ELIVITJ-X(2)-C-x-C.

NAME: Sodium:alanine symporter family signature.

20 CONSENSUS: G-G-x-EGAD(2)-ELIVMD-F-W-M-W-ELIVMD-x-ESTAVD-ELIVMFAD(2)-G.

NAME: Sodium:dicarboxylate symporter family signature l. CONSENSUS: P-x(0-l)-G-EDE3-x-ELIVMF3(2)-x-ELIVM3(2)-EKREQ3-ELIVM3(3)-x-P.

25 ELIVMI(3)-x-P.

NAME: Sodium:dicarboxylate sympostor family

NAME: Sodium:dicarboxylate symporter family signature 2.

CONSENSUS: P-x-G-x-ESTAJ-x-ENTJ-ELIVMCJ-D-G-ESTANJ-x-ELIVMJEFYJ-x(2)-ELIVMJ-x(2)-

30 CONSENSUS: ELIVMI-EFYI-ELII-ESAI-Q.

NAME: Sodium:galactoside symporter family signature.

CONSENSUS: D-x(3)-G-x(3)-EDN3-x(6-8)-G-EKH3-F-EKR3-P-EFYW3ELIVM3(2)-x-EGSTA3(2).

NAME: Sodium:neurotransmitter symporter family signature].
CONSENSUS: W-R-F-EGPI-Y-x(4)-N-G-G-G-x-EFYI.

NAME: Sodium:neurotransmitter symporter family signature 2.
40 CONSENSUS: Y-ELIVMFY3-x(2)-ESC3-ELIVMFY3-ESTQ3-x(2)-L-P-Wx(2)-C-x(4)-N-EGST3.

NAME: Sodium:solute symporter family signature 1.

CONSENSUS: EGSI-x(2)-ELIYI-x(3)-ELIVMFYWSTAGI(10)-ELIYIETAVI-x(2)-G-G-ELMFI-x-

45 ETAVJ-x(2)-G-G-ELMFJ-x-CONSENSUS: ESAPJ.

NAME: Sodium:solute symporter family signature 2.

CONSENSUS: EGASTI-ELIVMI-x(3)-EKRI-x(4)-G-A-x(2)-EGASI-

50 ELIVMGSJ-ELIVMWJ-ELIVMGATJ-G-CONSENSUS: x-ELIVMGJ.

NAME: Sodium:sulfate symporter family signature.

CONSENSUS: ESTACPI-S-x(2)-F-x(2)-P-ELIVMI-EGSAI-x(3)-N-x
55 ELIVMI-V.

NAME: glpT family of transporters signature. CONSENSUS: R-G-x(5)-W-N-x(2)-H-N-x-G-G.

NAME: Ammonium transporters signature.
CONSENSUS: D-FFYHST-A-C TCSC--(5)x-(5)EQAZI-(E)x-EVII-(5)x-ECZQI-Q-A-EZWYII-D

ESAGI-ELIVMFI-x(3)-

40

CONZENZUZ: ELIVMFYWAD(2)-x-EGKD-x-R.

BCCT family of transporters signature. CONZENZUZ: CGSDNJ-W-T-CLIVMJ-x-CFYJ-W-x-W-W-

10 Flagellar motor protein motA family signature. CONSENSUS: A-ELMF3-x-EGAT3-T-ELIVF3-x-G-x-ELIVMF3-x(7)-P.

Formate and nitrite transporters signature 1. CONZENZUZ: ELIVMAI-ELIVMYI-x-G-EGSTAI-EDESI-L-EFII-ETNI-EGSI.

15 Formate and nitrite transporters signature 2. conzensus: EGAD-x(2)-ECAD-N-ELIVMFYWD(2)-V-C-ELVD-A.

Prokaryotic sulfate-binding proteins signature 1. K-x-ENGEKI-EGTI-G-EDQI-x-ELIVMI-x(3)-Q-S. 20 CONZENZUZ:

NAME: Prokaryotic sulfate-binding proteins signature 2. CONSENSUS: N-P-K-ESTI-S-G-x-A-R.

NAME: Sulfate transporters signature.
CONSENSUS: P-x-Y-MCCT-1-V MCTT-1-V -(E)x-(E)EY7MVIJJ-(4)x-(5)EDATZJ-Y-L-CZDJ-Y----EGSTAD(2)-S-EKRD.

Amino acid permeases signature.

30 ESTAGCI-G-EPAGI-x(2,3)-ELIVMFYWAI(2)-x-ELIVMFYWI-CONZENZUZ: x-ELIVMFUSTAGCE(2)-CONZENZUZ: -EATOMILID-(E) x-ETZMVILID-x-ELIVMZIJ-x-(3)-ELIMZIJ-EGAB-E-x(5)-EPSALD.

35 Aromatic amino acids permeases signature. CONZENZUZ: $I-G-\mathbb{E}GA\mathbb{J}-G-M-\mathbb{E}LF\mathbb{J}-\mathbb{E}SA\mathbb{J}-x-P-x(3)-\mathbb{E}SA\mathbb{J}-G-x(2)-F.$

Xanthine/uracil permeases family signature. CONSENSUS: ELIVMD-P-x-EPASIFD-V-ELIVMD-G-G-x(4)-ELIVMD-EFYD-EGSAI-x-ELIVMI-x(3)-G.

Anion exchangers family signature 1.
F-G-G-ELIVMJ(2)-EKRJ-D-ELIVMJ-ERKJ-R-R-Y. CONZENZUZ:

45 Anion exchangers family signature 2. EFID-L-I-S-L-I-F-I-Y-E-T-F-x-K-L. NAME: CONZENZUZ:

MIP family signature. :SUZNZENCO ETABLE TABLE
50 General diffusion Gram-negative porins signature. NAME: CONZENZUZ: LLIVMFYD-x(2)-G-x(2)-Y-x-F-x-K-x(2)-ESND-ESTAVD-**ELIVMFYW3-V.**

55 NAME: OmpA-like domain. CONZENZUZ: -(5)x-Eqt23J-c2dd-(5)x-EAdJ-cATJ-x-ET3J-x-CAMVIJJ CLFYDEJ-ENQSJ-x(2)-

CONSENSUS: CLID-CSGD-CQED-CKRQED-R-A-x(2)-CLVD-x(3)-CLIVMFD-x(4,5)-CLIVMD-x(4)-

CONZENSUS: ELIVMI-x(3)-ESGI-x-G.

5 NAME: Eukaryotic mitochondrial porin signature.

CONSENSUS: EYHJ-x(2)-D-ESPAJ-x-ESTAJ-x(3)-ETAGJ-EKRJ-ELIVMFJ-

-(4)x-EZNGJ-EATZNGJ

CONSENSUS: EGSTAND-ELIVMAD-x-ELIVMYD.

10 NAME: Insulin-like growth factor binding proteins signature.

CONSENSUS: $G-C-\mathbb{E}GS\mathbb{J}-C-C-x(2)-C-A-x(F)-C$

NAME: GPR1/FUN34/yaaH family signature.
CONSENSUS: N-P-EAVI-P-ELFI-G-L-x-EGSAI-F.

15

NAME: GNS1/SUR4 family signature.

CONSENSUS: L-x-F-L-H-x-Y-H-H.

NAME: 43 Kd postsynaptic protein signature.

20 CONSENSUS: G-Q-D-Q-T-K-Q-Q-I.

NAME: Actins signature 1.

CONSENSUS: EFYD-ELIVD-G-EDED-E-A-Q-x-ERKQD(2)-G.

25 NAME: Actins signature 2.

CONSENSUS: W-EIVJ-ESTAJ-ERKJ-x-EDEJ-Y-EDNEJ-EDEJ.

NAME: Actins and actin-related proteins signature.

-CDATZON-V-UPHGAMY-X-ELIVND-T-E-EGAPQD-X-ELIVHFYUHQD-N-EPSTAQD-

30 \times (2)-N-EKRI.

NAME: Annexins repeated domain signature.

CONSENSUS: ETGJ-ESTVJ-x(B)-ELIVMFJ-x(2)-R-x(3)-EDEQNHJ-x(7)-

CIFY3-x(7)-CLIVMF3-

35 CONSENSUS: x(3)-ELIVMFD-x(11)-ELIVMFD-x(2)-ELIVMFD.

NAME: Caveolins signature.
CONSENSUS: F-E-D-V-I-A-E-P.

40 NAME: Clathrin light chain signature 1.

CONSENSUS: F-L-A-Q-Q-E-S.

NAME: Clathrin light chain signature 2.

CONSENSUS: EKRJ-D-x-S-EKRJ-ELIVMJ-EKRJ-x-ELIVMJ(3)-x-L-K.

45

NAME: Clusterin signature 1.

CONSENSUS: C-K-P-C-L-K-x-T-C.

NAME: Clusterin signature 2.

50 CONSENSUS: C-L-ERKI-M-ERKI-x-EEQI-C-EEDI-K-C.

NAME: Connexins signature 1.

CONSENSUS: C-EDNI-T-x-Q-P-G-C-x(2)-V-C-Y-D.

55 NAME: Connexins signature 2.

CONSENSUS: C-x(3,4)-P-C-x(3)-ELIVMJ-EDENJ-C-EFYJ-ELIVMJ-ESAJ-

EKRI-P.

WO 01/98454

CT/IB01/02050

Crystallins beta and gamma 'Greek key' motif NAME:

signature.

ELIVMFYWAD-x-{DEHRKSTP}-EFYD-EQHKYD-x(3)~EFYD-x-**CONZENZUZ:**

G-x(4)-ELIVMFCSTI.

5

NAME: Dynamin family signature.

CONSENSUS: L-P-ERKU-G-ESTNU-EGNU-ELIVMU-V-T-R.

NAME: Dynein light chain type 1 signature.

10 CONZENZUZ: H-x-I-x-G-EKRJ-x-F-EGAJ-S-x-V-ESTJ-EHYJ-E.

FtsZ protein signature 1.

CONSENSUS: N-ESTD-D-x-Q-x-L-x(16,18)-G-x-G-EATVD-G-EGSAND-x-

P-x(2)-G.

15

NAME: FtsZ protein signature 2.

EDNHKRJ-ELIVMFJ-x-ELIVMFJ(2)-EVZTACJ-EZTACJ-G-x-G-CONSENSUS:

EGKB-G-T-G-ESTD-G-

EGSARD-ESTAD-P-ELIVMFTD-ELIVMFD-ESGAVD. **CONSENSUS:**

20

NAME: Fungal hydrophobins signature.

CONSENSUS: EGNI-EDNQPSAI-x-C-EGSTANKI-EGSTADNQI-ESTNQII-

EPTIVI-x-C-C-EDENQKPSTI.

25 NAME: Intermediate filaments signature.

CONZENZUZ: CIVD-x-CTACID-Y-CRKHD-x-CLMD-L-CDED.

NAME: Involucrin signature.

CONSENSUS: <M-S-EQH3-Q-x-T-ELV3-P-V-T-ELV3.</p>

30

NAME: Kinesin motor domain signature.

CONZENZUZ: EGSAD-EKRHPSTQVMD-ELIVMFD-x-ELIVMFD-EIVCD-D-L-

EAHD-G-ESAND-E -

Kinesin motor domain profile. 35 NAME:

NAME: Kinesin light chain repeat.

EDEGRI-A-L-x(3)-EGEGI-x(3)-G-x-EDNSI-x-P-x-V-A-CONSENSUS:

-EZAJ-J-x-N-(E)x

40 **CONZENZUZ:** x(5)-EQR3-x-EKR3-EFY3-x(2)-EAV3-x(4)-EHKNQ3.

NAME: Myelin basic protein signature.

V-V-H-F-F-K-N. CONZENZUZ:

45 NAME: Myelin PO protein signature.

S-EKRJ-S-x-K-EAGJ-x-ESAJ-E-K-K-ESTAJ-K. **CONZENZUZ:**

NAME: Myelin proteolipid protein signature 1.

CONZENSUS: G-EMVD-A-L-F-C-G-C-G-H.

50

NAME: Myelin proteolipid protein signature 2.

CONSENSUS:

Α.

55 NAME: Neuromodulin (GAP-43) signature 1.

<M-L-C-C-ELIVMI-R-R. CONZENZUZ:

Neuromodulin (GAP-43) signature 2. NAME:

CONSENSUS: S-F-R-G-H-I-x-R-K-K-ELIVMI.

NAME: Osteopontin signature.

CONZENZUZ: EKQD-x-ETAD-x(2)-EGAD-S-S-E-E-K.

5 NAME: Peripherin / rom-l signature.

CONSENSUS: D-EGSJ-V-P-F-ESTJ-C-C-N-P-x-S-P-R-P-C.

NAME: Profilin signature.

10 CONSENSUS: <x(D¬Ъ)-EATAB-x(D¬Ъ)-W-EDRDHI-x-EYII-x-EDEQI-

NAME: Surfactant associated polypeptide SP-C palmitoylation

sites.

CONSENSUS: I-P-C-C-P-V.

15

NAME: Synapsins signature 1.

CONSENSUS: L-R-R-R-L-S-D-S.

NAME: Synapsins signature 2.

20 CONSENSUS: G-H-A-H-S-G-M-G-K-V-K-

NAME: Synaptobrevin signature.

CONSENSUS: N-ELIVMJ-EDANSJ-EKLJ-V-x-EDEQJ-R-x(2)-EKRJ-ELIVMJ-

ESTDEJ-x-ELIVMJ-x-EDEJ-

25 CONSENSUS: EKRI-ETAI-EDEI.

NAME: Synaptophysin / synaptoporin signature.

CONSENSUS: L-S-V-EDED-C-x-N-K-T-

30 NAME: Tropomyosins signature.

CONSENSUS: L-K-E-A-E-x-R-A-E-

NAME: Tubulin subunits alphan betan and gamma signature.

· SUZNEZNO : SUZNEZNO : SUZNEZNO :

35 ... NAME: Tubulin-beta mRNA autoregulation signal.

CONSENSUS: <M-R-EDEJ-EILJ.

NAME: Tau and MAP proteins tubulin-binding domain signature.

40 CONSENSUS: $G-S-x(2)-N-x(2)-H-x-\mathbb{E}PA\mathbb{J}-\mathbb{E}AG\mathbb{J}-G(2)$.

NAME: Neuraxin and MAPlB proteins repeated region signature.

CONSENSUS: ESTAGDNI-Y-x-Y-E-x(2)-EDEI-EKRI-ESTAGCII.

45 NAME: F-actin capping protein alpha subunit signature 1.

CONSENSUS: V-H-EFYI(2)-E-D-G-N-V.

NAME: F-actin capping protein alpha subunit signature 2.

CONSENSUS: F-K-EAEI-L-R-R-x-L-P.

NAME: F-actin capping protein beta subunit signature.

CONSENSUS: C-D-Y-N-R-D.

NAME: Vinculin family talin-binding region signature.

55 CONSENSUS: EKRI-x-ELIVMFI-x(3)-ELIVMAI-x(2)-ELIVMI-x(b)-R-Q-

Q-E-L.

50

NAME: Vinculin repeated domain signature.

CONSENSUS: ELIVMI-x-EQAI-A-x(2)-W-EILI-x-EDNI-P.

NAME: Amyloidogenic glycoprotein extracellular domain

signature.

CONSENSUS: G-EVTJ-E-EFYJ-V-C-C-P.

NAME: Amyloidogenic glycoprotein intracellular domain

signature.

CONSENSUS: G-Y-E-N-P-T-Y-EKRI.

10

40

NAME: Cadherins extracellular repeated domain signature.
CONSENSUS: ELIVI-x-ELIVI-x-D-x-N-D-ENHI-x-P.

NAME: Insect cuticle proteins signature.

15 CONSENSUS: G-x(7)-ENBURN-B-x(b)-Y-x-k-ENRUN-B-(2,3)-G-EYU-x-EQADI-

NAME: Gas vesicles protein GVPa signature 1.

CONSENSUS: CLIVMD-x-CDED-CLIVMFYTD-CLIVMD-CDED-x-CLIVMD(2)-

20 $\mathbb{C}DKRJ(2)-G-x-\mathbb{C}LIVMJ(2)$.

NAME: Gas vesicles protein GVPa signature 2.
CONSENSUS: R-ELIVAJ(3)-A-EGSJ-ELIVAJ-x-T-x(3)-Y-EAGJ.

25 NAME: Gas vesicles protein GVPc repeated domain signature. CONSENSUS: F-L-x(2)-T-x(3)-R-x(3)-A-x(2)-Q-x(3)-L-x(2)-F.

NAME: Bacterial microcompartiments proteins signature.

CONSENSUS: D-x(O₁L)-M-x-K-ESAGI(2)-x-EIVI-x-ELIVMI-ELIVMAI-

-EGG923-EG93-(4)x-E293 08
-EA93 : ZUZN3ZNO)

NAME: Flagella basal body rod proteins signature.
CONSENSUS: EGTARYQI-x(P)-ELIVMYSTAI(2)-EGTAI-ESTADENI-N-

35 ELIVMI-ESANI-N-x-ESADNFRI-CONSENSUS: ESTVI.

NAME: Flagella transport protein fliP family signature l. CONSENSUS: EPAB-A-EFYD-x-ELIVTD-ESTHD-EEQD-ELID-x(2)-EGAD-F-EKREQD-EIMD-G-ELIFD.

NAME: Flagella transport protein fliP family signature 2. CONSENSUS: P-ELIVMFB-K-ELIVMFB(5)-x-ELIVMAB-EDNGSB-G-W.

45 NAME: Plant viruses icosahedral capsid proteins 'S' region signature.

CONSENSUS: EFYWD-x-EPSTAD-x(7)-G-x-ELIVMD-x-ELIVMD-x-EFYWD-x(2)-D-x(5)-P.

NAME: Neurotransmitter-gated ion-channels signature.

55 CONSENSUS: C-x-ELIVMFQI-x-ELIVMFJ-x(2)-EFYI-P-x-D-x(3)-C.

NAME: ATP P2X receptors signature.

CONSENSUS: G-G-x-ELIVMI-G-ELIVMI-x-EIVI-x-W-x-C-EDNI-L-D-x(5)-C-x-P-x-Y-x-F.

NAME: G-protein coupled receptors signature.

5 CONSENSUS: EGSTALIVMFYWCI-EGSTANCPDEI-(EDPKRH)-x(2)-

ELIVMNQGAD-x(2)-ELIVMFTD-

CONSENSUS: EGSTANCI-ELIVMFYWSTACI-EDENHI-R-EFYWCSHI-x(2)-ELIVMI.

10 NAME: G-protein coupled receptors family 2 signature 1.

CONSENSUS: C-x(3)-EFYWLIVI-D-x(3,4)-C-EFWI-x(2)-ESTAGVI
x(8,9)-C-EPFI.

NAME: G-protein coupled receptors family 2 signature 2.

15 CONSENSUS: Q-G-ELMFCAD-ELIVMFTD-ELIVD-x-ELIVFSTD-ELIFDEVFYHD-C-ELFYD-x-N-x(2)-V.

NAME: G-protein coupled receptors family 3 signature 1.
CONSENSUS: [LV]-x-N-[LIVM](2)-x-L-F-x-I-[PA]-Q-[LIVM]-[STA]-

20 x-ESTAD(3)-ESTAND.

NAME: G-protein coupled receptors family 3 signature 2.

CONSENSUS: C-C-EFYWD-x-C-x(2)-C-x(4)-EFYWD-x(2,4)-EDND-x(2)ESTAHD-C-x(2)-C.

NAME: G-protein coupled receptors family 3 signature 3-CONSENSUS: F-N-E-ESTAD-K-x-I-ESTAGD-F-ESTD-M,

NAME: Visual pigments (opsins) retinal binding site.
30 CONSENSUS: ELIVMWACD-EPGACD-x(3)-ESACD-K-ESTALIMRD-EGSACPNVD-ESTACPD-x(2)-EDENFD-

CONSENSUS: EAPJ-x(2)-EIYJ.

NAME: Bacterial rhodopsins signature 1.
CONSENSUS: R-Y-x-EDTI-W-x-ELIVMFI-ESTI-T-P-ELIVMI(3).

EFY3.

35

40

NAME: Receptor tyrosine kinase class II signature-CONSENSUS: EDNI-ELIVI-Y-x(3)-Y-Y-R.

NAME: Receptor tyrosine kinase class III signature-45 CONSENSUS: G-x-H-x-N-ELIVMI-V-N-L-L-G-A-C-T-

NAME: Receptor tyrosine kinase class V signature L.
CONSENSUS: F-x-EDN3-x-EGAWJ-EGA3-C-ELIVM3-ESA3-ELIVM3(2)ESA3-ELV3-EKRHQ3-ELIVA3-

50 CONSENSUS: x(3)-EKRI-C-EPSAWI.

NAME: Receptor tyrosine kinase class V signature 2.

CONSENSUS: C-x(2)-EDED-G-EDEQD-W-x(2-3)-EPAQD-ELIVMTD-EGTD-x-C-x-C-x(2)-G-EHFYD-

55 CONSENSUS: EEQI.

NAME: Growth factor and cytokines receptors family signature 1.

WO 01/98454

РСТ/ІВ01/02050

CONSENSUS: C-ELVFYRD-x(7,8)-ESTIVDND-C-x-W.

NAME: Growth factor and cytokines receptors family signature

2.

5 CONSENSUS: ESTGLI-x-W-ESGI-x-W-S.

NAME: TNFR/NGFR family cysteine-rich region signature-CONSENSUS: C-x(4,b)-EFYHD-x(5,D)-C-x(0,2)-C-x(2,3)-C-

x(7,11)-C-x(4,6)-EDNEQSKP1-

10 CONSENSUS: x(5)-C.

NAME: TNFR/NGFR family cysteine-rich region domain-

NAME: Integrins alpha chain signature.

15 CONSENSUS: EFYWSII-ERKII-x-G-F-F-x-R.

NAME: Integrins beta chain cysteine-rich domain signature.

20 NAME: Natriuretic peptides receptors signature-CONSENSUS: G-P-x-C-x-Y-x-A-A-x-V-x-R-x(3)-H-W.

NAME: Photosynthetic reaction center proteins signature.

-H-x-EJDA23-(LL)x-EDA23-(S)x-H-x-9-(P)x-EHN3 :ZUZN3ZNOJ

25 ESAGI(2).

NAME: Antenna complexes alpha subunits signature.

-(Erd)x-EHMAT23-EBAT283

30 CONSENSUS: ESTNB-W-ELIVMFYWD.

NAME: Antenna complexes beta subunits signature.

CONSENSUS: $\mathbb{C}EQJ=x(4)-H=x(5)-\mathbb{C}QJ=x(3)-\mathbb{C}QJ=x($

[AV3-H-x(7)-P.

NAME: Photosystem I psaA and psaB proteins signature.

CONSENSUS: C-D-G-P-G-R-G-G-T-C.

NAME: Photosystem I psaG and psaK proteins signature.

40 CONSENSUS: G-F-x-ELIVMI-x-EDEAI-x(2)-EGAI-x-EGTAI-ESAI-x-G-H-

x-ELIVMI-EGAI.

NAME: Phytochrome chromophore attachment site signature.

NAME: Phytochrome chromophore attachment site domain

profile.

NAME: Speract receptor repeated domain signature.

50 CONSENSUS: G-x(5)-G-x(2)-E-x(6)-W-G-x(2)-C-x(3)-EYWJ-x(6)-C-

x(3)~G.

NAME: TonB-dependent receptor proteins signature 1.

CONSENSUS: <x(lo-lls)-EDENFI-ESTI-ELIVMFI-ELIVSTEQI-V-x-

55 EAGPI-ESTANEQPKI.

NAME: TonB-dependent receptor proteins signature 2.

CONSENSUS: ELYATIONATED-x-EQUATOUT : SUZNATOR TO SUZNA

5 NAME: Transmembrane 4 family signature.

CONSENSUS: G-x(3)-ELIVMFD-x(2)-EGSAD-ELIVMFD(2)-G-C-x-EGAD-ESTAD-x(2)-ECD-x(2)
CONSENSUS: ECWND-ELIVMD(2).

NAME: Bacterial chemotaxis sensory transducers signature.

CONSENSUS: R-T-E-EEQI-Q-x(2)-ESAI-ELIVMI-x-EEQI-T-A-A-S-M-E-Q-L-T-A-T-V.

NAME: ER lumen protein retaining receptor signature 1.

CONSENSUS: G-I-S-x-EKRI-x-Q-x-L-EFYI-x-ELIVI(2)-F-x(2)-R-Y-

NAME: ER lumen protein retaining receptor signature 2-CONSENSUS: L-E-ESAI-V-A-I-ELMI-P- ϱ -L.

20 NAME: Ephrins signature.

CONSENSUS: EKRQI-ELFI-ECSTI-x-K-EIFI-Q-x-EFYI-ESTI-EPAI-x(3)G-x-E-F-x(5)-EFYI(2)CONSENSUS: x(2)-ESAI.

25 NAME: Granulins signature. CONSENSUS: C-x-D-x(2)-H-C-C-P-x(4)-C.

30

NAME:

NAME: HBGF/FGF family signature.

CONSENSUS: G-x-L-x-ESTAGP3-x(6,7)-EDE3-C-x-EFM3-x-E-x(6)-Y.

NAME: PTN/MK heparin-binding protein family signature 1.
CONSENSUS: S-EDE3-C-x-EDE3-W-x-W-x(2)-C-x-P-x-ESN3-x-D-C-GELIVMA3-G-x-R-E-G.

35 NAME: PTN/MK heparin-binding protein family signature 2. CONSENSUS: C-EKRI-ELIVMI-P-C-N-W-K-K-x-F-G-A-EDEI-C-K-Y-x-F-EEQI-x-W-G-x-C.

NAME: Nerve growth factor family signature.
40 CONSENSUS: G-C-EKRI-G-ELIVI-EDEI-x(3)-EYWI-x-S-x-C.

NAME: Platelet-derived growth factor (PDGF) family signature.

CONSENSUS: P-EPSI-C-V-x(3)-R-C-EGSTAI-G-C-C.

NAME: Small cytokines (intercrine/chemokine) C-x-C subfamily signature.

CONSENSUS: C-x-C-ELIVMD-x(5,b)-ELIVMFYD-x(2)-ERKSEQD-x-ELIVMD-x(3)-ELIVMD-x(5)-

50 CONSENSUS: ESAGI-x(2)-C-x(3)-EEQI-ELIVMI(2)-x(9,10)-C-L-EDNI.

signature.
CONSENSUS: C-C-ELIFYTD-x(5,b)-ELID-x(4)-ELIVMFD-x(2)-EFYWD-

Small cytokines (intercrine/chemokine) C-C subfamily

NAME: TGF-beta family signature.

CONSENSUS: ELIVM3-x(2)-P-x(2)-EFY3-x(4)-C-x-G-x-C.

NAME: TNF family signature.

CONSENSUS: ELVJ-x-ELIVMJ-x(3)-G-ELIVMFJ-Y-ELIVMFJ(2)-x(2)-

5 EQEKHLI-ELIVMGTI-x-

CONSENSUS: ELIVMFY1.

NAME: TNF family profile.

10 NAME: Unt-1 family signature.

CONSENSUS: C-K-C-H-G-ELIVMTD-S-G-x-C.

NAME: Interferon alpha, beta and delta family signature.

CONSENSUS: EFYHI-EFYI-x-EGNRCI-ELIVMI-x(2)-EFYI-L-x(7)-ECYI-

15 A-W.

NAME: Granulocyte-macrophage colony-stimulating factor

signature.

CONSENSUS: C-P-ELPI-T-x-E-ESTI-x-C.

20

NAME: Interleukin-1 signature-

x(7)-ELIVMD.

25 NAME: Interleukin-2 signature.

CONSENSUS: T-E-ULFI-x(2)-L-x-C-L-x(2)-E-L

NAME: Interleukins -4 and -13 signature.

CONSENSUS: L-x-E-ELIVMJ(2)-x(4,5)-ELIVMJ-ETLJ-x(5,7)-C-x(4)-

30 ELVAD-x-EDNDD-ELIVMAD.

NAME: Interleukin-b / G-CSF / MGF signature.
CONSENSUS: C-x(9)-C-x(b)-G-L-x(2)-EFYI-x(3)-L.

35 NAME: Interleukin-7 and -9 signature.
CONSENSUS: N-x-ELADJ-ESCTJ-F-L-K-x-L-L.

NAME: Interleukin-10 family signature.

40

NAME: LIF / OSM family signature.

NAME: Macrophage migration inhibitory factor family

45 signature.

CONSENSUS: EDEI-P-C-A-x(3)-ELIVMI-x-S-I-G-x-ELIVMI-G.

NAME: Adipokinetic hormone family signature.

CONSENSUS: Q-ELV3-ENT3-EFY3-EST3-x(2)-W.

50

NAME: Bombesin-like peptides family signature.

CONSENSUS: W-A-x-G-ESHJ-ELFJ-M.

NAME: Calcitonin / CGRP / IAPP family signature.

-ETYYJ-CAMVI-J-T-EAZI-(L, C)x-ENTZI-ENGDAZI-C(E)

x(3)-ELYF3-

NAME: Corticotropin-releasing factor family signature.

CONSENSUS: EPQ3-x-ELIVM3-S-ELIVM3-x(2)-EPST3-ELIVMF3-x-ELIVM3-L-R-x(2)-ELIVM3-

NAME: Crustacean CHH/MIH/GIH neurohormones family signature5 CONSENSUS: C-EDENKI-D-C-x-N-ELIVI-EFYI-R-x(7)-C-EKRI-x(2)-C.

NAME: Erythropoietin / thrombopoeitin signature.
CONSENSUS: P-x(4)-C-D-x-R-ELIVM1(2)-x-EKR1-x(14)-C-

10 NAME: Granins signature 1.
CONSENSUS: EDEJ-ESNJ-L-ESANJ-x(2)-EDEJ-x-E-L.

NAME: Granins signature 2.

CONSENSUS: C-ELIVMJ(2)-E-ELIVMJ(2)-S-EDNJ-ESTAJ-L-x-K-x-S-

15 x(3)-ELIVMD-ESTAD-x-E-C.

30

40

NAME: Galanin signature.
CONSENSUS: G-W-T-L-N-S-A-G-Y-L-L-G-P-H.

20 NAME: Gastrin / cholecystokinin family signature-CONSENSUS: Y-x(0-1)-EGDI-EWHI-M-EDRI-F.

NAME: Glucagon / GIP / secretin / VIP family signature.
CONSENSUS: EYHI-ESTAIVGDI-EDEQI-EAGFI-ELIVMSTEI-EFYLRI-x-

25 CDENSTAKI-EDENSTAICONSENSUS: ELIVATYGI-×(9)-EKREQLI-EKRDENQLI-ELVFYWGI-ELIVQI.

NAME: Glycoprotein hormones alpha chain signature L. CONSENSUS: C-x-G-C-C-EFYD-S-R-A-EFYD-P-T-P.

NAME: Glycoprotein hormones alpha chain signature 2.
CONSENSUS: N-H-T-x-C-x-C-x-T-C-x(2)-H-K.

NAME: Glycoprotein hormones beta chain signature 1.
35 CONSENSUS: ___ C-ESTAGMI-G-EHFYLI-C-x-ESTI.

NAME: Glycoprotein hormones beta chain signature 2.

CONSENSUS: EPAI-V-A-x(2)-C-x-C-x(2)-C-x(4)-ESTDI-EDEYI-C-x(6-8)-EPGSTAVMI-x(2)-C.

NAME: Gonadotropin-releasing hormones signature.
CONSENSUS: Q-H-EFYWJ-S-x(4)-P-G.

NAME: Insulin family signature.

45 CONSENSUS: C-C-{P}-x(2)-C-ESTDNEKPID-x(3)-ELIVMFSD-x(3)-C.

NAME: Natriuretic peptides signature.

CONSENSUS: C-F-G-x(3)-D-R-I-x(3)-S-x(2)-G-C.

NAME: Neurohypophysial hormones signature.
CONSENSUS: C-ELIFYI(2)-x-N-ECSI-P-x-G.

NAME: Neuromedin U signature.
CONSENSUS: F-ELIVMFI-F-R-P-R-N.

NAME: Endogenous opioids neuropeptides precursors signature.

CONSENSUS: C-x(3)-C-x(2)-EKRH3-x(b,7)-ELIF3-EDN3-x(3)
C-x-ELIVM3-EEQ3-C-

CONZENZUZ: $\mathbb{E}Q\mathbb{J}-x(B)-U-x(2)-C$.

NAME: Pancreatic hormone family signature.

CONSENSUS: CFYD-x-(3)-CMVIJD-x(2)-Y-x(3)-CLIVMFYD-x-R-x-R-

5 EYFI.

NAME: Parathyroid hormone family signature.

CONSENSUS: V-S-E-x-Q-x(2)-H-x(2)-G.

10 NAME: Pyrokinins signature. CONSENSUS: F-EGSTV3-P-R-L-EG>1.

NAME: Somatotropin, prolactin and related hormones signature 1.

15 CONSENSUS: C-x-ESTD-x(2)-ELIVMFYD-x-ELIVMSTAD-p-x(5)-ETALIVD-x(7)-ELIVMFYD-x(6)-

CONZENZUZ: ELIVMFYJ-x(2)-EZTAJ-W.

NAME: Somatotropin, prolactin and related hormones signature 20 2.

CONSENSITS:

CONSENSUS: C-ELIVMYYI-x(2)-D-ELIVMYYSTAI-x(5)-ELIVMYYI-x(2)-ELIVMYYII-x(2)-C-

NAME: Tachykinin family signature-25 CONSENSUS: F-EIVFYI-G-ELMI-M-EG>I.

NAME: Thymosin beta-4 family signature. CONSENSUS: K-L-K-T-E-T-Q-E-K-N.

30 NAME: Urotensin II signature. CONSENSUS: C-F-W-K-Y-C.

NAME: Cecropin family signature.

CONSENSUS: W-x(0,2)-EKDND-x(2)-K-EKRED-ELID-E-ERKND.

NAME: Mammalian defensins signature.

CONSENSUS: C-x-C-x(3,5)-C-x(7)-G-x-C-x(9)-C-C.

NAME: Arthropod defensins signature.

40 CONSENSUS: C-x(2,3)-EHNJ-C-x(3,4)-EGRJ-x(2)-G-G-x-C-x(4,7)-C-x-C-

NAME: Cathelicidins signature 1.

CONSENSUS: Y-x-EEDJ-x-V-x-ERQJ-A-ELIVMAJ-EDQGJ-x-ELIVMYJ-N-

45 [EQ].

×(3)-C.

NAME: Cathelicidins signature 2.

CONSENSUS: F-x-ELIVMJ-K-E-T-x-C-x(JD)-C-x-F-EKRJ-EKEJ.

50 NAME: Endothelin family signature.
CONSENSUS: C-x-C-x(4)-D-x(2)-C-x(2)-EFYI-C.

NAME: Plant thionins signature.
CONSENSUS: C-C-x(5)-R-x(2)-EFYD-x(2)-C.

NAME: Gamma-thionins family signature.

CONSENSUS: EKRI-x-C-x(3)-ESVI-x(2)-EFYWHI-x-EGFI-x-C-x(5)-C-

NAME: Snake toxins signature.

CONSENSUS: G-C-x(1,3)-C-P-x(8,10)-C-C-x(2)-EPDENI.

5 NAME: Myotoxins signature.
CONSENSUS: K-x-C-H-x-K-x(2)-H-C-x(2)-K-x(3)-C-x(3)-K-x(2)-Cx(2)-ERK3-x-K-C-C-K-K.

NAME: Scorpion short toxins signature.

10 CONSENSUS: C-x(3)-C-x(6,9)-EGASI-K-C-EIMQTI-x(3)-C-x-C.

NAME: Heat-stable enterotoxins signature.

15 NAME: Aerolysin type toxins signature. CONSENSUS: EKTI-x(2)-N-W-x(2)-T-EDNI-T.

CONZENZUZ:

NAME: Shiga/ricin ribosomal inactivating toxins active site signature.

20 CONSENSUS: CLIVMAI-x-CLIVMSTAI(2)-x-E-ESAGVI-ESTALI-R-EFYIERKNQSI-x-CLIVMI-EEQSICONSENSUS: ...x(2)-ELIVMI.

C-C-x(2)-C-C-x-P-A-C-x-G-C.

NAME: Channel forming colicins signature25 CONSENSUS: T-x(2)-W-x-P-ELIVMFY1(3)-x(2)-E-

NAME: Staphylococcal enterotoxin/Streptococcal pyrogenic exotoxin signature 1.
CONSENSUS: Y-G-G-ELIVI-T-x(4)-N.

35 NAME: Staphyloccocal enterotoxin/Streptococcal pyrogenic exotoxin signature 2.

CONSENSUS: K-x(2)-ELIVI-x(4)-ELIVI-D-x(3)-R-x(2)-L-x(5)-ELIVI-Y.

40 NAME: Thiol-activated cytolysins signature.
CONSENSUS: ERKB-E-C-T-G-L-x-W-E-W-U-ERKB.

NAME: Membrane attack complex components / perforin signature.

45 CONSENSUS: Y-x(b)-EFYJ-G-T-H-EFYJ.

NAME: Pancreatic trypsin inhibitor (Kunitz) family signature.

CONSENSUS: F-x(3)-G-C-x(b)-EFYI-x(5)-C.

NAME: Bowman-Birk serine protease inhibitors family signature.

CONSENSUS: C-x(5,6)-EDENQKRHSTAI-C-EPASTDHI-EPASTDKI-EASTDVI-C-ENDKSI-EDEKRHSTAI-C.

NAME: Kazal serine protease inhibitors family signature.
CONSENSUS: C-x(7)-C-x(b)-Y-x(3)-C-x(2-3)-C.

NAME: Soybean trypsin inhibitor (Kunitz) protease inhibitors family signature.

CONSENSUS: ELIVMD-x-D-x-EEDNTYD-EDGD-ERKHDENQD-x-ELIVMD-x(5)-

Y-x-ELIVMI.

NAME: Serpins signature.

----EYRMYND-EZQHYND-EZQHYND-X-EYRMVIJ-X-EYRMVIJ-X-EZHZHZZHOZHOZH

ELIVMFYD-ELIVMFYCD-x-

CONSENSUS: ELIVMFAHI.

10

NAME: Potato inhibitor I family signature.

CONSENSUS:

EFYWJ-P-EEQHJ-ELIVJ(2)-G-x(2)-ESTAGVJ-x(2)-A.

NAME: Squash family of serine protease inhibitors signature.

15 CONSENSUS: C-P-x(5)-C-x(2)-D-x-D-C-x(3)-C-x-C.

NAME: Streptomyces subtilisin-type inhibitors signature.

C-x-P-x(2,3)-G-x-H-P-x(4)-A-C-EATDD-x-L.

20 NAME: Cysteine proteases inhibitors signature.

CONSENSUS: EGSTECKRVI-Q-ELIVTI-EVAFI-ESAGCI-G-x-ELIVMNKI
x(2)-ELIVMFYI-x-ELIVMFYAI
CONSENSUS: EDENCKRHSIVI.

25 NAME: Tissue inhibitors of metalloproteinases signature. CONSENSUS: C-x-C-x-P-x-H-P-Q-x-A-F-C.

NAME: Cereal trypsin/alpha-amylase inhibitors family signature.

30 CONSENSUS: C-x(4)-ESAGDI-x(4)-ESALI-ELFI-x(2)-C-EHI-x-ELIVMFYI(2)-x(3,4)-C.

NAME: Alpha-2-macroglobulin family thiolester region signature.

35 CONSENSUS: EPGJ-x-EGSJ-C-EGAJ-E-EEQJ-x-ELIVMJ.

NAME: Disintegrins signature.

CONSENSUS: C-x(2)-G-x-C-C-x-ENQRSD-C-x-EFMD-x(b)-C-ERKD.

40 NAME: Lambdoid phages regulatory protein CIII signature. CONSENSUS: E-S-x-L-x-R-x(2)-EKRI-x-L-x(4)-EKRI(2)-x(2)-EDEI-x-L.

NAME: Chaperonins cpnb0 signature.

45 CONSENSUS: A-EASI-x-EDEQI-E-x(4)-G-G-EGAI.

NAME: Chaperonins cpnlO signature.

CONSENSUS: CLIVMFYD-x-P-CILTD-x-CDEND-CKRD-CLIVMFAD(3)-

EKREQU-x(B,9)-ESGU-x-

50 CONSENSUS: ELIVMFYI(3).

NAME: Chaperonins TCP-1 signature 1.

CONSENSUS: CRKELI-CSTI-x-CLMFYI-G-P-x-CGSAI-x-x-K-CLIVMFI(2).

55 NAME: Chaperonins TCP-1 signature 2.
CONSENSUS: ELIVM1-ETS1-ENK1-D-EGA1-EAVNHK1-ETAV1-ELIVM1(2)x(2)-ELIVM1-x-ELIVM1-xCONSENSUS: ESNH1-EPQH1.

PCT/IB01/02050 WO 01/98454

Chaperonins TCP-1 signature 3-NAME:

CONZENZUZ: Q-EDEKD-x-x-ELIVMGTAD-EGAD-D-G-T.

5 NAME: Heat shock hsp20 proteins family profile.

NAME: Heat shock hsp?D proteins family signature 1.

CONSENSUS: EIV3-D-L-G-T-EST3-x-ESC3.

10 NAME: Heat shock hsp?O proteins family signature 2.

-CT2AJ-C2DJ-CH2DJ-G-R27MVIJJ-EVQJ-EV7MVIJJ-C7MVIJJ-CONZENZUZ:

-EMVIJU-ETZU-(E)x

CONZENSUS: **LLIVMFCI**.

Heat shock hsp70 proteins family signature 3. 15 NAME: : SUZNAZNO ELIVMYD-x-ELIVMFD-x-G-G-x-ESTD-x-ELIVMD-P-x-

ELIVMD-x-EDERKRSTAD.

Heat shock hsp90 proteins family signature.

CONSENSUS: Y-x-ENGHJ-K-EDEJ-EIVAJ-F-L-R-EEDJ-20

Chaperonins clpA/B signature 1.

CONZENZUZ: D-EAID-ESGAD-N-ELIVMFI(2)-K-EPTD-x-L-x(2)-G.

25 NAMF: Chaperonins clpA/B signature 2.

R-ELIVMFYJ-D-x-C-ELYMVIJ-3-X-E-EKRQJ-x-ESTAJ-x-CONSENSUS:

ESTAD-EKRD-ELIVMD-x-G-: SUZNAZNO . EATZI

Nt-dnaJ domain signature. 30

CONZENZUZ: EFY3-x(2)-ELIVMAJ-x(3)-EFYWHNTJ-EDENQSAJ-x-L-x-

EDNU-x(3)-EKRU-x(2)-EFYIU-

NAME: dnaJ domain profile.

CXXCXGXG dnaJ domain signature.

CONSENSUS: C-EDEGSTHKRD-x-C-x-G-x-EGKD-EAGSDMD-x(2)-EGSNKRD-

x(4-6)-C-x(2-3)-C-x-G-x-G.

40

NAME: grpE protein signature.
CONSENSUS: EFLI-FDNII-FDNII-EFLD-EDND-EPHEAD-x(2)-EHMD-x-A-ELIVMTND-x(16,20)-

 $G-\mathbb{E}[Y]]-x(3)-\mathbb{E}[DEG]]-x(2)-$

CONSENSUS: ELIVMI-ERII-x-ESAI-x-V-x-EIVI-

45 NAME: Bacterial type II secretion system protein C

signature.

35

CONZENZUS: P-x(b)-F-x(4)-L-x(3)-D-ELIVMD-A-ELIVMD-x-ELIVMD-N-

x-ELIVMI-x-L.

NAME: 50 Bacterial type II secretion system protein D

signature.

CONSENSUS: EGRI-EDERKGI-ESTVMI-ELIVMAI(3)-EGAI-G-ELIVMFYI-

 $\times (JJ) - ELIVMJ - P -$

CONSENSUS: CLIVMFYWGSD-CLIVMFD-CGSAED-x-CLIVMD-P-

CLIVMFYW3(2)-x(2)-ELV3-F-55

> NAME: Bacterial type II secretion system protein E

signature.

CUNION : CLIVMILAR -x (2) -P-D-x-ELIVMI(3) -G-E-ELIVMILAR -D.

NAME: Bacterial type II secretion system protein F

signature.

CONSENSUS: EKAGI-EFINWAI-x-(5)-EZAIAI-EFINAI-x-ETAI-b-x(5)-

LLIVMD-x(B)-ESTAGVD-x(b)-

CONSENSUS: [LMY]-x(3)-[LIVMF](2)-P.

NAME: Bacterial type II secretion system protein N

10 signature.

CONSENSUS: G-T-L-W-x-G-x(11)-L-x(4)-W.

NAME: Bacterial export FHIPEP family signature.

15 P-G-K-Q-M-EGSAI-I-D-

CONSENSUS: EGSAI-D..

NAME: Protein secA signatures.

CONSENSUS: EIVJ-x-EIVJ-ESAJ-T-ENQJ-M-A-G-R-G-x-D-I-x-L.

20 NAME: Protein secY signature 1.

CONSENSUS: EGST3-ELIVM3-x-P-

ELIVMFYD(2)-x-EASD-EGSTQD-

CONSENSUS: CLIVMFAT(2)-Q-ELIVMFAT(2).

25

NAME: Protein secY signature 2.

CONSENSUS: ELIVMYIJJ(2)-x-EDEJ-x-ELIVMFJ-ESTNJ-x(2)-G-

CLIVMFD-CGSTD-CNSTD-G-x-CGSTD-

CONSENSUS: ELIVMFI(3).

30

NAME: Protein secE/secbl-gamma signature.

CONSENSUS: LLIVHYJ-x(2)-LADBNAGJ-x(4)-LATMVIJ-x-EATMVIJ-x(2)-

 $\mathbb{C}KWJ-P-x(3)-\mathbb{C}SEQJ-x(7)-$

CONSENSUS: ELIVID-ELIVGAD-ELIVFGASTD.

35

NAME: Gram-negative pili assembly chaperone signature.

CONSENSUS: ELIVMFYJ-EAPNJ-x-EDNSJ-EKREQJ-E-ESTRJ-ELIVMARJ-x-

EFYWT3-x-ENC3-ELIVM3-

CONZENZUZ: x(2)-ELIVMJ-P-EPAZJ.

40

NAME: Fimbrial biogenesis outer membrane usher protein

signature.

CONSENSUS: EVLI-EPASQI-EPASI-G-EPADI-EFYI-x-ELII-EDNQSTAPI-

CDNHJ-CLIVMFYJ.

45

SRP54-type proteins GTP-binding domain signature.

CONSENSUS: P-ELIVMJ-x-EFYLJ-ELIVMJ-x-EGSJ-x-EGSJ-EEQJ-x(4)-

ELIVMF3.

NAME:

50 NAME: Cytochrome c oxidase assembly factor COX10/ctaB/cvoE

signature.

CONSENSUS: $\mathbb{L}EDJ-x-D-x(2)-M-x-R-T-x(2)-R-x(4)-G$.

NAME: Cyclin-dependent kinases regulatory subunits signature

55 1.

CONSENSUS: Y~S~x~EKRJ~Y~x~EDEJ(2)-x~EFYJ~E~Y~R~H~V~x~ELVJ~

EPTI-EKRPI.

NAME: Cyclin-dependent kinases regulatory subunits signature

2.

CONSENSUS: H-x-P-E-x-H-EIVI-L-L-F-EKRI.

5 NAME: Pentaxin family signature. CONSENSUS: H-x-C-x-ESTI-W-x-ESTI-

NAME: Immunoglobulins and major histocompatibility complex proteins signature.

10 CONSENSUS: EFYJ-x-C-x-EVAJ-x-H.

NAME: Prion protein signature 1.

CONSENSUS: A-G-A-A-A-G-A-V-V-G-G-L-G-G-Y.

15 NAME: Prion protein signature 2.

CONSENSUS: E-x-EED3-x-K-ELIVM3(2)-x-EKR3-ELIVM3(2)-x-EQE3-MC-x(2)-Q-Y.

NAME: Cyclins signature.

20 CONSENSUS: R-x(2)-ELIVMSAI-x(2)-EFYWSI-ELIVMI-x(8)-ELIVMFCI-x(4)-ELIVMFYAI-x(2)CONSENSUS: ESTAGCI-ELIVMFYQI-x-ELIVMFYCI-ELIVMFYI-D-ERKHI-ELIVMFYWI.

25 NAME: Proliferating cell nuclear antigen signature l.

CONSENSUS: EGAD-ELIVMFD-x-ELIVMAD-x-ESAVD-ELIVMD-D-x-ENSAEDEHKRD-EVID-x-ELYDCONSENSUS: EVGAD-x-ELIVMD-x-ELIVMD-x(4)-F.

30 NAME: Proliferating cell nuclear antigen signature 2.

CONSENSUS: ERKAD-C-EDED-ERHD-x(3)-ELIVMFD-x(3)-ELIVMD-xESGAND-ELIVMFD-x-KCONSENSUS: ELIVMFD(2).

35 NAME: Actin-depolymerizing proteins signature.
CAYD-X-EAD-X-EAD-X-ELIVNTD-EKRD-X-EKRD-M-ELIVMD-EYADEAYD-X-ERMVID-CAYD-X-EKRD-X-EKRD-M-ELIVMD-EYADEAYD-X-EMAND-X-EXA

40 NAME: BCL2-like apoptosis inhibitors (spans part of BH3, BH) and BH2).

NAME: Apoptosis regulator, Bcl-2 family BHL domain signature.

45 CONSENSUS: ELYMED-EFTD-x-EGSDD-EGLD-x(L-2)-ENSD-EYWD-G-R-ELIVD-ELIVCD-EGTD-CONSENSUS: ELIVMFB(2)-x-F-EGSAED-EGSARYD.

NAME: Apoptosis regulator, Bcl-2 family BH2 domain

50 signatureCONSENSUS: W-ELIMB-x(3)-EGRD-G-EWQB-EDENSAVD-x-EFLGADELIVFTCD:

NAME: Apoptosis regulator, Bcl-2 family BH3 domain

55 signature.
CONSENSUS: ELIVATI-x(3)-L-EKARQI-x-EIVALI-G-D-EDESGI-ELIMFVIEDENSHQI-ELVSHRQICONSENSUS: ENSRI-

NAME: Apoptosis regulator, Bcl-2 family BH4 domain

signature.

CONSENSUS: EDSJ-ENTJ-R-EAEJ-ELIJ-V-x-EKDJ-EFYJ-ELIVJ-EGHSJ-Y-

5 K-L-ESRI-Q-ERKI-G-

CONZENZUZ: EHYI-x-ECWI.

NAME: Apoptosis regulator, Bcl-2 family BH4 domain profile.

10 NAME: Arrestins signature.

CONSENSUS: EFYD-R-Y-G-x-EDED(2)-x-EDED-ELIVMD(2)-G-ELIVMD-x-F-x-ERKD-EDEQD-ELIVMD.

NAME: AAA-protein family signature.

15 CONSENSUS: ELIVATI-x-ELIVATI-ELIVAFI-x-EGATACI-ESTI-ENSI-

x(4)-ELIVMD-D-x-A-ELIFAD-

CONSENSUS: x-R.

NAME: Ubiquitin domain signature.

20 CONSENSUS: K-x(2)-ELIVMD-x-EDESAKD-x(3)-ELIVMD-EPAD-x(3)-Q-x-

ELIVMI-ELIVMCI-

CONSENSUS: ELIVMFYD-x-G-x(4)-EDED.

NAME: Ubiquitin domain profile.

25

NAME: ADP-ribosylation factors family signature.

ELIVNJ-x(2)-EGSAD-ELIVNFJ-x-

CONSENSUS: EMKI-ELIVMI.

30

NAME: GTP-binding nuclear protein ran signature.

CONSENSUS: D-T-A-G-Q-E-K-ELFJ-G-G-L-R-EDEJ-G-Y-Y.

NAME: SARL family signature.

35 CONSENSUS: R-x-ELIVMI-E-V-F-M-C-S-ELIVMI(2)-x-EKRQI-x-G-Y-x-

E-EAGI-EFII-x-W-ELIVMI-

CONSENSUS: x-Q-Y.

NAME: Band 7 protein family signature.

40 CONSENSUS: R-x(2)-ELIVI-ESANI-x(b)-ELIVI-D-x(2)-T-x(2)-U-G-

CLIVI-CKRHI-CLIVI-x-

CONSENSUS: EKR3-ELIV3-E-ELIV3-EKR3.

NAME: Trp-Asp (WD) repeats signature.

45 CONSENSUS: ELIVMSTACI-ELIVMFYWSTAGCI-ELIMSTAGI-ELIVMSTAGCI-

x(2)-EDNI-x(2)-

CONSENSE: ELIVMUSTACI-x-ELIVMFSTAGI-W-EDENJ-ELIVMFSTAGCNJ.

NAME: G-protein gamma subunit profile.

50

NAME: Ras GTPase-activating proteins signature.

CSUZNAZA - CLIVMFIJ-CLIVMFYJ-R-CLIVMFYJ-R-CLIVMFYJ(2)-

EGACND-P-EAVD-ELIVD(2)-

CONSENSUS: ESGAND-P.

55

NAME: Ras GTPase-activating proteins profile.

NAME: Guanine-nucleotide dissociation stimulators CDC24 family signature.

CONSENSUS: L-x(2)-ELIVMT-L-x(2)-P-ELIVMT-x(2)-ELIVMT-x-

EKRSI-x(2)-L-x-ELIVMI-x-

5 CONSENSUS: EDEQI-ELIVMI-x(3)-ESTI.

NAME: Guanine-nucleotide dissociation stimulators CDC25

family signature.

CONSENSUS: EGAPI-ECTI-V-P-EFYI-x(4)-ELIVMFYI-x-EDNI-ELIVMI.

10

NAME: MARCKS family signature 1.
CONSENSUS: G-Q-E-N-G-H-V-EKRI.

NAME: MARCKS family phosphorylation site domain.

15 CONSENSUS: E-T-P-K(5)-x(0,1)-F-S-F-K-K-x-F-K-L-S-G-x-S-F-K-K-RI-ENSI-EKRI-K-E

NAME: Stathmin family signature 1.

CONZENZUZ: P-EKGJ-EKRJ(2)-EDEJ-x-Z-L-EEGJ-E.

20

NAME: Stathmin family signature 2-CONSENSUS: A-E-K-R-E-H-E-EKRI-E-V.

NAME: GTP-binding elongation factors signature.

25 CONSENSUS: D-EKRSTGANQFYWD-x(3)-E-EKRAQD-x-ERKQDD-EGCD-EIVMKD-ESTD-EIVD-x(2)-

CONSENSUS: EGSTACKRNQI.

NAME: Elongation factor 1 beta/beta//delta chain signature

30 1.

CONSENSUS: EDED-EDEGD-EDED(2)-ELIVMFD-D-L-F-G.

NAME: Elongation factor 1 beta/beta'/delta chain signature

2.

35 CONSENSUS: V-Q-S-x-D-ELIVM3-x-A-EFWM3-ENQ3-K-ELIVM3.

NAME: Elongation factor 1 gamma chain profile.

NAME: Elongation factor Ts signature 1.

40 CONSENSUS: L-R-x(2)-T-EGDQJ-x-EGSJ-ELIVMFJ-x(0,1)-EDRNKACJ-x-K-EKRNEQSJ-EAVJ-L-

NAME: Elongation factor Ts signature 2.

CONSENSUS: É-ELIVMI-N-ESCUI-EQÈI-T-D-F-V-ESAI-EKRNI.

45

NAME: Elongation factor P signature.

CONSENSUS: K-x-A-x(4)-G-x(2)-ELIVJ-x-V-P-x(2)-ELIVJ-x(2)-G.

NAME: Eukaryotic initiation factor lA signature.

50 CONSENSUS: EIMI-x-G-x-EGSI-EKRHI-x(4)-ECLI-x-D-G-x(2)-R-x(2)ERHI-I-x-G.

NAME: Eukaryotic initiation factor 4E signature.

CONSENSUS: EDEII-EIFYII-x(2)-F-EKRII-x(2)-ELIVMII-x-P-x-W-E-EDVII-

55 x(5)-G-G-EKRI-W.

NAME: Eukaryotic initiation factor 5A hypusine signature. CONSENSUS: EPTI-G-K-H-G-x-A-K.

NAME: Initiation factor 2 signature.

CONSENSUS: G-x-ELIVM3-x(2)-L-EKRJ-EKRHNS3-x-K-x(5)-ELIVM3-

x(2)-G-x-EDENI-C-G.

NAME: Initiation factor 3 signature.

CONSENSUS: [KR]-[LIVM](2)-[DN]-[FY]-[GSN]-[KR]-[LIVMFYS]-x[FY]-[DEQT]-x(2)-[KR].

10 NAME: Translation initiation factor SUIL signature. CONSENSUS: ELIVMJ-EEQJ-ELIVMJ-Q-G-EDENJ-EKHQJ-EKRVJ.

NAME: Prokaryotic-type class I peptide chain release factors signature.

15 CONSENSUS: EARI-ESTAI-x-G-x-G-Q-EHNGCSI-V-N-x(3)-ESTI-A-EIVI.

20

NAME: Calponin family repeat.

CONSENSUS: ELIVID-x-ELSJ-Q-EMASJ-G-ESTYJ-ENTJ-EKRQJ-x(2)ESTNJ-Q-x-G-x(3,4)-G.

25 NAME: CAP protein signature 1-CONSENSUS: ELIVMI(2)-x-R-L-EDEI-x(4)-R-L-E-

NAME: CAP protein signature 2.

CONSENSUS: D-ELIVMFY3-x-E-x-EPA3-x-P-E-Q-ELIVMFY3-K.

30

NAME: Calreticulin family signature 1.

CONSENSUS: EKRHNI-x-EDEQNI-EDEQNKI-x(3)-C-G-G-EAGI-EFYIELIVMI-EKNI-ELIVMFYI(2).

35 NAME: Calreticulin family signature 2.
CONSENSUS: ELIVMI(2)-F-G-P-D-x-C-EAGI

NAME: Calreticulin family repeated motif signature.

CONSENSUS: EIVJ-x-D-x-EDENSTJ-x(2)-K-P-EDEHJ-D-W-EDENJ.

40

NAME: Calsequestrin signature 2.

45 CONSENSUS: IDEI-L-E-D-W-ELIVMI-E-D-V-L-x-G-x-ELIVMI-N-T-E-D-D-D-D-

50 EFYD-x-EESD-EFYVCD-x(2)-CONSENSUS: ELIVMFSD-ELIVMFD.

NAME: Hemolysin-type calcium-binding region signature. CONSENSUS: D-x-ELID-x(4)-G-x-D-x-ELID-x-G-G-x(3)-D.

55

CONSENSUS: ELIVMFYWD(2)-x-ELIVMFYWD(3).

NAME: P-II protein urydylation site.

CONSENSUS: Y-EKRJ-G-EASJ-EAEJ-Y.

NAME: P-II protein C-terminal region signature.

CONSENSUS: CZJJ-x(3)-G-EYQJ-G-EKRJ-EIVJ-EFWJ-ELTZJ-x(2)-ELTZJ-x(2)-

5

40

50

55

10 NAME: 14-3-3 proteins signature 1.

CONSENSUS: R-N-L-ELIVI-S-EVGI-EGAI-Y-EKNI-N-EIVAI.

NAME: 14-3-3 proteins signature 2.

CONSENSUS: Y-K-EDEJ-S-T-L-I-EIMJ-Q-L-ELFJ-ERHCJ-D-N-ELFJ-T-

15 ELSI-W-ETANI-ESADI.

NAME: ATPLGB / PLM / MATB family signature.

CONSENSUS: EDNSI-x-F-x-Y-D-x(2)-ESTI-ELIVMI-ERQI-x(2)-G.

20 NAME: BTG1 family signature 1.

CONSENSUS: Y-x(2)-EHPJ-W-CYJ-EAPJ-E-x-P-x-K-G-x-EGAJ-EYJ-R-C-EIVJ-EHJ-EIVJ.

NAME: BTG1 family signature 2.

25 CONSENSUS: ELVI-P-x-EDEI-ELMI-ESTI-ELIVMI-W-EIVI-D-P-x-E-V-ESCI-x-ERQI-x-G-E.

NAME: Cullin family signature.

CONZENZUZ: ELIV3-K-x(2)-ELIV3-x(2)-L-I-EDEQ3-EKRHNQ3-x-Y-

30 CLIVMJ-x-R-x(6,7)-CFYJ-x-CLAZJ-x-Y : ZUZNJZNO)

NAME: Cullin family profile.

35 NAME: Enhancer of rudimentary signature.

CONZENZUZ: Y-D-I-ESAJ-x-L-EYJ-x-F-EIVJ-D-x(3)-D-ELIVJ-S.

NAME: GLO protein signature 1.

L-C-C-x-EKRI-C-x(4)-EDEI-x-N-x(4)-C-x-C-R-V-P.

NAME: GlO protein signature 2.

CONSENSUS: C-x-H-C-G-C-EKRHJ-G-C-ESAJ.

NAME: Glucokinase regulatory protein family signature.

45 CONSENSUS: G-EPAJ-E-x-ELIVJ-ESTAJ-G-S-ESTJ-R-ELIVMJ-K-

 $\mathbb{L}STGAJ(3) - \times (2) - K$.

NAME: GTP1/OBG family signature.

CONSENSUS: D-ELIVMJ-P-G-ELIVMJ(2)-EDEYJ-EGNJ-A-x(2)-G-x-G.

NAME: HIT family signature.

CONSENSUS: ENGAD-x-(4)-EGAVD-x-EQFD-x-ELIVMD-x-H-ELIVMYD-H-

CLIVMFTI-H-CLIVMFI(2)-

CONSENSUS: EPSGAI.

NAME: Caseins alpha/beta signature.
CONSENSUS: C-L-ELV3-A-x-A-ELVF3-A.

NAME: Clathrin adaptor complexes medium chain signature L-CONSCNUS: EIVTD-EGSPD-W-R-x(2,3)-EGADD-x(2)-EHYD-x(2)-N-x-ELIVND-

CONSENSUS: ELIVATI-E.

NAME: Clathrin adaptor complexes medium chain signature 2. CONSENSUS: ELIVD-x-F-I-P-P-x-G-x-ELIVMFYD-x-L-x(2)-Y.

NAME: Clathrin adaptor complexes small chain signature.

10 CONSENSUS: CLIVMJ(2)-Y-EKRJ-x(4)-L-Y-F.

NAME: Ependymins signature 1.

CONSENSUS: F-E-E-G-x-ELIVMFJ-Y-EEDJ-I-D-x(2)-N-EQEJ-S-C-

ERKHI(2).

15

NAME: Ependymins signature 2.

CANDENCE EQUIDATE CATE OF CATE OF CAMPILLAND STATE OF CATE OF

20 NAME: Syntaxin / epimorphin family signature.

CONSENSUS: CRQ3-x(3)-CLNVIJ-(2)-CLNVIJ-(2)-CLIVNI-

__x-EDEVMJ-ELIVMJ-x(2)-

-SUZNAZZIONE (E)x-ETVIJ-x(2)x-EZ7J-EMVIJJ-x(2)-Q-

EGADEQU-x(2)-ELIVMU-EDNQTU-x-

25 CONSENSUS: ELIVMFI-EDESVI-x(2)-ELIVMI.

NAME: Extracellular proteins SCP/Tpx-l/Ag5/PR-l/Sc7

signature L. CONSENSUS:

EGDERJ-H-EFYWHJ-T-Q-ELIVMJ(2)-W-x(2)-ESTNJ.

30

NAME: Extracellular proteins SCP/Tpx-l/Ag5/PR-l/Sc7 signature 2.

CONSENSUS: ELIVMFYH3-ELIVMFYJ-x-C-ENQRHSJ-Y-x-EPARH3-x-EGL3-N-ELIVMFYH3-ELIVMFYH3-ELIVMFYH3-x-EGL3-

35

NAME: Fetuin family signature 1.

CONSENSUS: C-x(5b)-C-x(10)-C-x(13)-C-x(17-18)-C-x(13)-C-x(2)-

C-x(58)-C-x(10,11)-

CONZENZUZ: C-x(JD-JZ)-C-x(JP-ZS)-C.

40

NAME: Fetuin family signature 2.

CONSENSUS: L-E-T-x-C-H-x-L-D-P-T-P

CONSENSUS: L-E-T-x-C-H-x-L-D-P-T-P.

NAME: Legume lectins beta-chain signature.
45 CONSENSUS: ELIVI-ESTAGI-V-EDEQVI-EFLII-D-ESTI.

NAME: Vertebrate galactoside-binding lectin signature.
CONSENSUS: W-EGEKI-x-EEQI-x-EKREI-x(3,6)-EPCTFI-ELIVMFI-

ENGEGSKVI-x-EGHI-x(3)-

CONSENSUS: EDENKHSJ-ELIVMFCJ.

55 NAME: Lysosome-associated membrane glycoproteins duplicated domain signature.

CONSENSUS: ESTAI-C-ELIVMI-ELIVMFYWI-A-x-ELIVMFYWI-x(3)-

ILIVMFYWI-x(3)-Y.

NAME: LAMP glycoproteins transmembrane and cytoplasmic domain signature.

CONSENSUS: C-x(3)-D-x(3-4)-CLIVMJ(3)-P-CLIVMJ-x-CLIVMJ-G-

5 x(2)-ELIVMコ-x-G-ELIVMコ(2)-CONSENSUS: x-ELIVMコ(4)-A-EFYコ-x-ELIVMコ-x(2)-EKRコ-ERHコ-x(ユっ2)-ESTAGコ(2)-Y-EEQコ・

NAME: Glycophorin A signature.

10 CONSENSUS: I-I-x-EGACD-V-M-A-G-ELIVMD(2).

NAME: PMP-22 / EMP / MP20 family signature 1-

CONSENSUS: ELIVIDATE (4)-EZAJ-T-x(2)-EDNKSJ-x-W-x-EAJ-D-ELIVJ-W-

x(≥)-C-

40

NAME: PMP-22 / EMP / MP2D family signature 2.

CARD-EAVD-x-T-CVD-L-S-x-ELID-x(4)-EGSAD-EAVD-X-ELID-x(4)-EGSAD-ELAVD-X-ELID-X(4)-EGSAD-X-ELID-X(4)-EGSAD-X-ELAV

20 NAME: 0xysterol-binding protein family signature.
CONSENSUS: E-EKQI-x-S-H-EHRI-P-P-x-ESTACFI-A.

NAME: Yeast PIR proteins repeats signature.

CONSENSUS: S-Q-EIV3-ESTGNH3-D-G-Q-ELIV3-Q-EAIV3-ESTA3.

NAME: Seminal vesicle protein I repeats signature.
CONSENSUS: EIVMD-x-G-Q-D-x-V-K-x(5)-EKND-G-x(3)-ESTLVD.

NAME: Seminal vesicle protein II repeats signature.
30 CONSENSUS: EGSAD-Q-x-K-S-EFYD-x-Q-x-K-ESAD.

NAME: Serum amyloid A proteins signature.

CONSENSUS: A-R-G-N-Y-EEDI-A-x-EQKRI-R-G-x-G-G-x-W-A.

35 NAME: Spermadhesins family signature 1. CONSENSUS: C-G-x(2)-ELID-x(4)-G-x-I-x(9)-C-x-W-T.

NAME: Spermadhesins family signature 2.
CONSENSUS: C-x-K-E-x-ELIVMI-E-ELIVMI-x-EDEI-x(3)-EGSI-x(5)-K-x-C.

NAME: Stress-induced proteins SRP1/TIP1 family signature.
CONSENSUS: P-W-Y-EST1(2)-R-L.

45 NAME: Glypicans signature.

CONSENSUS: C-x(2)-C-x-G-ELIVMI-x(4)-P-C-x(2)-EFYI-C-x(2)-ELIVMI-x(2)-G-C.

NAME: Syndecans signature.
50 CONSENSUS: EFYD-R-EIMD-EKRD-K(2)-D-E-G-S-Y.

NAME: Tissue factor signature.
CONSENSUS: W-K-x-K-C-x(2)-T-x-EDENI-T-E-C-D-ELIVMI-T-D-E.

55 NAME: Translationally controlled tumor protein signature 1.
CONSENSUS: EIAB-G-EGASD-N-EPAD-S-A-E-EGDED-EPAGED-x(0,1)EDEGD-x-EDEND-x(2)-EDED.

NAME: Translationally controlled tumor protein signature 2. CONSENSUS: EFLI-EFYI-EIVTI-G-E-x-EMAI-x(2.5)-EDENI-EGASI-x-

CLVD-CAVD-x(3)-CFYD-CKRD-

CONZENZUZ: [DE]-

5

NAME: Tub family signature 1.

CONSENSUS: F-EKHQD-G-R-V-ESTD-x-A-S-V-K-N-F-Q.

NAME: Tub family signature 2.

10 CONSENSUS: A-F-EAGD-I-ESACD-ELIVMD-ESTD-S-F-x-EGSTD-K-x-A-C-

F.

NAME: HCP repeats signature.

CONSENSUS: H-R-H-R-G-H-x(2)-EDEJ(7).

15

NAME: Bacterial ice-nucleation proteins octamer repeat.

CONSENSUS: A-G-Y-G-S-T-x-T.

NAME: Cell cycle proteins ftsW / rodA / spoVE signature.

20 CONSENSUS: ENVI-x(5)-EGTRI-ELIVMAI-x-P-EPTLIVMI-x-G-ELIVMI-

-[AZYJ-Z-(5)[W7MVIJJ-(E)x

CONSENSUS: G-G-ESTND-ESAD.

NAME: Enterobacterial virulence outer membrane protein

25 signature 1.

CONSENSUS: G-ELIVMFYD-N-ELIVMD-K-Y-R-Y-E.

NAME: Enterobacterial virulence outer membrane protein

signature 2.

30 CONSENSUS: EFYWD-x(2)-G-x-G-Y-EKRD-F>.

NAME: Hydrogenases expression/synthesis hypA family

signature.

CONSENSUS: F-ECSAJ-EYJ-EDEJ-ELVAJ(2)-x(3)-ETJ-ELVVJ-

35 x(1b)-C-x(2)-C-x(12-15)-

CONSENSUS: C-P-x-C.

NAME: Hydrogenases expression/synthesis hupF/hypC family

signature.

40 CONSENSUS: <M-C-ELIVI-EGAI-ELIVI-P-x-EQKRI-ELIVI.

NAME: Staphylocoagulase repeat signature-

 $-D-\mathbb{E}N\mathbb{C}\mathbb{I}$ = $-D-\mathbb{E}$

. O-Y-G.

45

NAME: 11-5 plant seed storage proteins signature.

CONSENSUS: N-G-x-CDEJ(2)-x-CLIVMFJ-C-CSTJ-x(JJ-J2)-CPAGJ-D.

NAME: Dehydrins signature 1.

50 CONSENSUS: S(5)-EDEI-x-EDEI-G-x(1,2)-G-x(0,1)-EKRI(4).

NAME: Dehydrins signature 2-

CONSENSUS: EKRI-ELIMI-K-EDEI-K-ELIMI-P-G.

55 NAME: Germin family signature.

CONSENSUS: G-x(4)-H-x-H-P-x-A-x-E-ELIVMJ.

NAME: Oleosins signature.

-d-1-ECASI-EMVIJ-(2)x-EDAI-(2)x-ETZI-EDAI

CLIVMFJ(4)-F-S-P-ELIVMJ(3)-

CONSENSUS: P-A.

5 NAME: Small hydrophilic plant seed proteins signature. CONSENSUS: G-EEQI-T-V-V-P-G-G-T.

NAME: Pathogenesis-related proteins BetvI family signature.

CONSENSUS: G-x(2)-ELIVMFI-x(4)-E-x(2)-ECSTAENI-x(8,9)-EGNDI-

10 $G-\mathbb{E}GS\mathbb{J}-\mathbb{E}CS\mathbb{J}-\times(2)-K-\times(4)-$

CONSENSUS: [FY].

NAME: Pollen proteins Ole e I family signature.

CONSENSUS: EEQI-G-x-V-Y-C-D-T-C-R.

15

NAME: Thaumatin family signature.

CONSENSUS: $G-x-\mathbb{E}GF\mathbb{I}-x-C-x-T-\mathbb{E}GA\mathbb{I}-D-C-x(\mathbb{I}_1\mathbb{Z})-G-x(\mathbb{Z}_1\mathbb{Z})-C$.

NAME: Mrp family signature.

20 CONSENSUS: W-x(2)-ELIVMJ-D-ELIVMYJ(4)-D-x-P-P-G-T-EGSJ-D.

NAME: Glucose inhibited division protein A family signature

L.

CONSENSUS: EGSJ-P-x-Y-C-P-S-ELIVMJ-E-x-K-ELIVMJ-x-EKRJ-F.

NAME: Glucose inhibited division protein A family signature

CONSENSUS: A-G-Q-x-ENTI-G-x(2)-G-Y-x-E-ESAGI(3)-EQSI-G-

ELIVMI(2)-A-G-ELIVMTI-N-A.

30

NAME: NOL1/NOP2/sun family signature-

CONSENSUS: EFVI-D-EKRAI-ELIVMAI-L-x-D-EAVI-P-C-ESTI-EGAI.

NAME: PETIL2 family signature.

35 CONSENSUS: LDN3-x-LDN3-x-x(3)-P-L-LLIV3-E-LLIV3-x-LT3-x-P.

NAME: Protein smpB signature.

CONSENSUS: ETA3-G-ELIVM3-x-L-x-G-x-E-ELIVM3-EKQ3-ESA3-ELIVM3.

40 NAME: Hypothetical cof family signature 1.

CONSENSUS: ELIVATION - ELIVATI

ELVYD-ESTANLMD.

NAME: Hypothetical cof family signature 2.

45 CONSENSUS: ELIVMFCD-G-D-EGSANQD-x-N-D-x(3)-ELIMFYD-x(2)-EAVD-

 \times (2) \times (2) \times (2) \times

CONSENSUS: [LMP]-x(2)-[GAS].

NAME: RIO1/ZK632.3/MJ0444 family signature.

50 CONZENSUS: ELIVMI-V-H-EGAI-D-L-S-E-EFYI-N-x-ELIVMI.

NAME: SUA5/yci0/yrdC family signature.

CONSENSUS: ELIVATAD(3)-ELIVATYCD-EPGD-T-EDED-ESTAD-x-EFYD-

- EZAI-ELIVMI-EGZI-

55

NAME: Uncharacterized protein family UPF0001 signature.

CONSENSUS: EFWJ-H-EFMJ-EIVJ-G-x-ELIVJ-Q-x-ENKRJ-K-x(3)-ELIVJ.

NAME: Uncharacterized protein family UPF0003 signature.

CONSUSCIENCE: G-x-V-x(2)-ELIVI-x(3)-ESAI-x(b)-D-x(3)-ELIVII(3)-P-N-x(2)-ELIVIFI(2)
CONSUSCIENCE: x(5)-N.

NAME: Uncharacterized protein family UPFOOO4 signature.

CONSENSUS: ELIVMD-x-ELIVMTD-x(2)-G-C-x(3)-C-ESTAND-EFYD-C-xELIVMD-x(4)-G.

10 NAME: Uncharacterized protein family UPFOOOS signature.
CONSENSUS: G-ELIVMI(2)-ESAI-x(5,8)-G-x(2)-ELIVMI-G-P-x-Lx(4)-ESAI-x(4,6)CONSENSUS: ELIVMI(2)-x(2)-A-x(3)-T-A-ELIVMI(2)-F-

15 NAME: Uncharacterized protein family UPF0006 signature 1. CONSENSUS: ELIVMFYI(2)-D-ESTAI-H-x-H-ELIVMFI-EDNI.

NAME: Uncharacterized protein family UPFOOOL signature 2.
CONSENSUS: P-ELIVMD-x-ELIVMD-H-x-R-x-ETAD-x-EDED.

NAME: Uncharacterized protein family UPFODOL signature 3.
CONSENSUS: ELVSJ-ELIVAJ-x(2)-ELIVAJ-x(3)-L-ELIVAJELIVAJ-E-T-D-x-P.

25 NAME: Uncharacterized protein family UPFDB07 signature. CONSENSUS: V-L-EIVI-H-D-EGAI-A-R.

30

NAME: Uncharacterized protein family UPFOOll signature-CONSENSUS: S-D-A-G-x-P-x-ELIVI-ESNI-D-P-G.

NAME: Uncharacterized protein family UPFOOl2 signature.
CONSENSUS: USGRANUTE CONSENSUS: USGRA

NAME: Uncharacterized protein family UPFDD15 signature.

35 CONSENSUS:_ EDEJ-ELIVMFJ(3)-R-T-ESGJ-G-x(2)-R-x-S-x-EFYJ-ELIVMJ(2)-W-Q.

NAME: Uncharacterized protein family UPFOOL6 signature.
CONSENSUS: E-ELIVM3-G-D-K-T-F-ELIVMF3(2)-A.

40

NAME: Uncharacterized protein family UPFOOL? signature.

CONSENSUS: D-x(A)-EGNI-ELFYI-x(4)-EDETI-ELYI-Y-x(3)-ESTI
x(?)-EIVI-x(2)-EPSI-x
CONSENSUS: ELIVMI-x(E)-EDNI-D.

AS

NAME: Uncharacterized protein family UPFOOL9 signature.

CONSENSUS: L-P-V-EVTI-ENQLI-F-EATI-A-G-G-ELIVI-A-T-P-A-D-A-A-ELMI.

50 NAME: Uncharacterized protein family UPFOO2O signature-CONSENSUS: D-P-ELIVMFD-C-G-ESTD-G-x(3)-ELID-E.

NAME: Uncharacterized protein family UPFOD21 signature. CONSENSUS: C-K-x(2)-F-x(4)-E-x(22,23)-S-G-G-K-D.

NAME: Uncharacterized protein family UPFOO23 signature.
CONSENSUS: D-x-D-E-ELIVI-L-x(4)-V-F-x(3)-S-K-G-

NAME: Uncharacterized protein family UPF0024 signature. CONSENSUS: G-x-K-D-EKR3-x-A-ELV3-T-x-Q-x-ELIVF3-ESGC3.

NAME: Uncharacterized protein family UPFDO25 signature.

5 CONSENSUS: D-V-ELIVI-x(2)-G-H-ESTI-H-x(12)-ELIVMFI-N-P-G.

NAME: Uncharacterized protein family UPFOO27 signatureCONSENSUS: Q-ELIVM3-x-N-x-A-x-ELIVM3-P-x-I-x(b)-ELIVM3-P-D-xH-x-G-x-G-x(2)-EIV3-G.

10
NAME: Uncharacterized protein family UPFODZB signature.
CONSENSUS: EGAD-EGSD-G-EGAD-A-R-G-x-ESAD-H-x-G-x(9)-EIVD-x-EIVD-x-(2)-EGAD-G-x-S-

15
NAME: Uncharacterized protein family UPFOO29 signature.
CONSENSUS: G-x(2)-ELIVMI(2)-x(2)-ELIVMI-x(4)-ELIVMI-x(5)ELIVMI(2)-x-R-EFYWI(2)-GCONSENSUS: G-x(2)-ELIVMI-G.

NAME: Uncharacterized protein family UPFOO3D signature.
CONSENSUS: EGAB-L-I-ELIVI-P-G-G-E-S-T-ESTAI.

NAME: Uncharacterized protein family UPFOO31 signature 125 CONSENSUS: ESAVI-EIVWI-ELVAI-ELIVI-G-EPNSI-G-L-EGPI-xEDENGTI-

NAME: Uncharacterized protein family UPFOO31 signature 2-CONSENSUS: EGAD-G-x-G-D-ETVD-ELTD-ESTAD-G-x-ELIVMD.

NAME: Uncharacterized protein family UPFOO32 signature.
CONSENSUS: Y-x(2)-F-ELIVMAD(2)-x-L-x(4)-G-x(2)-F-EEQDELIVMFD-P-ELIVMD.

35 NAME:—— Uncharacterized protein family-UPF0033 signature. CONSENSUS: L-EDNI-x(2)-ETAGI-x(2)-C-P-x-P-x-ELIVMI.

NAME: Uncharacterized protein family UPF0034 signature.
CONSENSUS: ELIVM3-EDNG3-ELIVM3-N-x-G-C-P-x(3)-ELIVM3-x(5)-G-ESAC3.

NAME: Uncharacterized protein family UPF0035 signature.
CONSENSUS: L-L-T-x-R-ESAI-x(3)-R-x(3)-G-x(3)-F-P-G-G.

45 NAME: Uncharacterized protein family UPFOO36 signature.

CONSENSUS: H-x-S-G-H-EGAI-x(3)-EDEI-x(3)-ELMI-x(5)-P-x(3)
ELIVMI-P-x-H-G-EDEI.

NAME: Uncharacterized protein family UPFDD38 signature.
50 CONSENSUS: G-x-ELI3-x-R-x(2)-L-x(4)-F-x(8)-ELIV3-x(5)-P-x-ELIV3-

NAME: Uncharacterized protein family UPFOO44 signature.

CONSENSUS: L-ESTI-x(3)-K-x(3)-EKRI-ESGAI-x-EGAI-H-x-L-x-PELIVI-x(2)-ELIVI-EGAI-

55 ELIVI-x(2)-ELIVI-EGAI-CONSENSUS: x(2)-G.

CONZENZUZ:

30

40

x-G-

NAME: Uncharacterized protein family UPF0047 signature.

CONSENSUS: S-X(2)-ELIVI-x-ELIVI-x(2)-G-x(4)-G-T-W-Q-x-ELIVI.

NAME: Uncharacterized protein family UPF0054 signature.

CONZENZUZ: H-EGZJ-x-L-H-L-ELIJ-G-EFYWJ-D-H.

NAME: Uncharacterized protein family UPFDD57 signature.

CONSENSUS: ELIVI-x-ESTAI-ELIVFI(3)-P-P-ELIVAI-EGAI-EIVI-x(4)EGKNI.

10 NAME: Hypothetical YERD57c/yjjV family signature.

CONSENSUS: P-EATJ-R-ESAJ-x-ELIVMYJ-x(2)-EAKJ-x-L-P-x(4)ELIVMJ-E.

NAME: Hypothetical hesB/yadR/yfhF family signature.

15 CONSENSUS: F-x-ELIVMFYJ-x-N-EPGJ-ENSKJ-x(4)-C-x-C-EGSJ-x-S-F.

NAME: Hypothetical yabo/yceC/sfhB family signature.
CONSENSUS: ENHYD-R-ELID-D-x(2)-T-ESTD-G-ELIVMAD-ELIVMFD(2)ELIVMFGD-ESGACD.

20

30

Deposit of Clones

Each clone has been transfected into separate bacterial cells (E. coli) in the composite deposit.

The clones are located and publically available from the Resource Center of the German Human Genome Project (Heubner Weg by 14059 Berlin, GERMANY), from which each clone comprising a particular polynucleotide is obtainable. The Resource Center library numbers are slightly different that those presented here, but may be readily obtained by the following key or with the assistance of Resource Center personnel.

The library name becomes a number: brain (hfbr2) becomes

564; kidney (hfkd2) becomes 566; mammary carcinoma (hmcfl)

becomes 727; testis (htes3) becomes 434; amygdala (hamy2) becomes

761; melanoma (hmel2) becomes 762 and uterus (hutel) becomes 586.

Next; the plate number is converted to two digits (e.g., "2"

becomes "02") and is moved behind the plate coordinate; and the

underscore is dropped. The following examples are helpful:

Listed Number Resource Center Number DKFZphamy2_l0hl7 DKFZp761H1710 DKFZphfbr2_78i21 DKFZp564I2178 45 DKFZphfkd2_3k1 DKFZp566K013 DKFZphmcfl_lc23 DKFZp727C231 DKFZhmel2_12j1 DKFZp7b2J0112 DKFZphtes3_16b5 DKFZp434B0516 DKFZphutel_17k7 DKFZp586K0717

10

15

20

25

35



The libraries were constructed using two commercially available vectors. The brain (hfbr2 designations) and kidney (hfkd2 designations) libraries utilize pAMP 1 from Life Technologies and are maintained in XL-2Blue (Strategene); the amyqdala (hamy2), testes (htes3) and melanoma (hmel2) libraries are constructed in pSPORTL, also from Life Technologies, and are maintained in DH10B (LifeTechnologies). In addition to the following techniques, consultation with the commercial literature available on these clones will make evident all of the housekeeping techniques needed to propagate and isolate the individual constructs. All inserts may be excised with a NotI/SalI digestion. Alternatively, universal primers, flanking the cloning region, may be used to amplify the inserts using PCR methods.

Bacterial cells containing a particular clone can be obtained from the composite deposit as follows:

An oligonucleotide probe or probes should be designed to the sequence that is known for that particular clone. This sequence can be derived from the sequences provided herein, or from a combination of those sequences. Methods of probe design are presented below.

Oligonucleotide probes may be labeled with -32P ATP (specific activity 6000 Ci/mmole) and T4 polynucleotide kinase using commonly employed techniques for labeling oligonucleotides. Other, non-radioactive labeling techniques can also be used. Unincorporated label typically is removed by gel filtration chromatography or other established methods. The amount of radioactivity incorporated into the probe can be quantified by 30 measurement in a scintillation counter. Preferably, specific activity of the resulting probe generally should be approximately 4X10b dmp/pmole-

The bacterial culture containing the pool of full-length clones should preferably be thawed and 100 l of the stock used to inoculate a sterile culture flask containing 25 ml of sterile L-broth containing ampicillin at 50 - 100 g/ml (for XL-2Blue strains 25 g/ml tetracycline should also be used). The culture should preferably be grown to saturation at 37°C., and the saturated culture should preferably be diluted in fresh L-broth.

Aliquots of these dilutions should preferably be plated to determine the dilution and volume which will yield approximately 5000 distinct and well-separated colonies on solid bacteriological media containing L-broth containing ampicillin at 100 g/ml (for XL-2Blue strains 25 g/ml tetracycline should also be used) and agar at 1.5% in a 150 mm petri dish when grown overnight at 37°C. Other known methods of obtaining distinct, well-separated colonies can also be employed.

Standard colony hybridization procedures should then be used to transfer the colonies to nitrocellulose filters and lyse-10 denature and bake them. The filter is then preferably incubated at 65°C. for 1 hour with gentle agitation in 6 x SSC (20 x stock is 175.3 g NaCl/liter, 88.2 g Na citrate/liter, adjusted to pH 7.D with NaOH) containing D.5% SDS, LOO g/ml of yeast RNA, and 15 10 mM EDTA (approximately 10 mL per 150 mm filter). Preferably, the probe is then added to the hybridization mix at a concentration greater than or equal to 1X10 dpm/mL. The filter is then preferably incubated at 65°C. with gentle agitation overnight. The filter is then preferably washed in 500 mL of 2 x 20 SSC/0-5% SDS at room temperature without agitation, preferably followed by 500 mL of 2 x SSC/0.1% SDS at room temperature with gentle shaking for 15 minutes. A third wash with 0.1 x SSC/0.5% SDS at 65°C. for 30 minutes to 1 hour is optional. The filter is then preferably dried and subjected to autoradiography for sufficient time to visualize the positives on the X-ray film-25 Other known hybridization methods can also be employed.

The positive colonies are picked, grown in culture, and plasmid DNA isolated using standard procedures. The clones can then be verified by restriction analysis, hybridization analysis, or DNA sequencing.

Alternatively, clones may be grown as described above, and PCR used to isolate the insert DNAs. Methods of PCR are described below and are otherwise well known.

ERROR SCREENING

30

35

The DNA sequences found herein derive from individual clones, which are publicly available, as noted above. Thus, the skilled artisan will recognize that any specific sequence disclosed herein

readily can be screened for errors by resequencing a particular fragment, in both directions (i.e., by sequencing both strands). Alternatively, error screening can be performed by amplifying and/or cloning any of the inventive DNAs, using for example RT-PCR, and sequencing the resulting amplified product. In the event that there is a sequencing error, reference should be made to the deposited clone as the correct sequence.

USES AND BIOLOGICAL ACTIVITIES OF THE INVENTIVE MOLECULES

The inventive molecules and their derivatives are susceptible to a wide variety of uses, based on functional and/or structural properties. The skilled worker will appreciate, based on the biological activities detailed below, and discussed with regard to the individual sequences herein, that the inventive molecules will find usefulness in numerous therapeutic and diagnostic applications.

The DNA molecules, especially the potassium salts thereofican be used as fertilizer supplements due to their high nitrogen and phosphorus contents. Since the DNAs are of defined length, they are also useful in gel electrophoresis as molecular weight markers. Due to their similarity with known molecules, certain of the DNA molecules and their variants and derivatives may be used in any number of different diagnostic procedures and therapeutic applications. They may also be used to make the encoded proteins.

The proteins themselves have many possible uses. They may be used as a nutritional supplement for humans, animals and even for laboratory use as, for example, medium for bacterial cultures. Moreover, since the proteins are of defined, known sizes, they may be used as molecular weight markers for gel electrophoresis and gel filtration. Because they are of defined sequences, they also have use in microsequencing and protein fingerprinting applications.

Expression Profiling Applications

10

15

20

25

30

35

Given their known tissue expression and functional associations, assemblages of the inventive proteins (or corresponding antibodies) and nucleic acids are particularly suited to expression profiling applications. Expression profiling generally entails constructing an array of indicators that signal

-484-

the presence of a particular RNA or protein expression product. Such arrays can be used to evaluate for example pharmacological effectiveness and toxicity. In particular expression profiles from such arrays can be generated from cells treated with known compounds having known properties and these profiles can be compared to profiles of unknowns to evaluate similarities and differences which can be correlated with efficacy or toxicity.

Additional uses of profiling include diagnosis, tracking development, and ascertaining signaling and metabolic pathways.

10

15

20

25

30

35

For examples of references describing profiling and its uses, see Farr et al., U.S. Patent 5,811,231 (1998); Seilhamer et al., U.S. Patent 5,840,484 (1998); Rine et al., U.S. Patent No. 5,777,888 (1998); WO 97/27317; WO 99/05323; WO 99/09218; and WO 99/14369. For a device for implementing such techniques, see Lipshutz et al., U.S. Patent No. 5,856,174 (1999) and Anderson et al., U.S. Patent No. 5,922,591 (1999).

In one embodiment, a subset of the inventive DNAs will be arrayed on a substrate, like a gene chip, a filter or a 96-well plate. Test samples containing cells are maintained in the presence of a label capable of incorporation into nascent mRNA. Samples are treated with test and control compounds, which will induce mRNA expression in the sample, resulting in incorporation of label. Whole mRNA is isolated and applied to the array such that it hybridizes with the DNAs contained therein. After washing, the amount of hybridization is quantified and a profile is generated. These steps are repeated with various control and test compounds, thereby generating a library of profiles, which can be used to ascertain the relationships relevant to pharmacological efficacy or toxicity.

The matrices used in such profiling, however, need not be limited to those utilizing DNAs. Rather, other nucleic acids, like RNAs and protein nucleic acids (PNAs), as well as the inventive proteins and antibodies corresponding to the inventive proteins may also be employed. Hence, for example, antibodies could form the array and the samples could be treated in order to label nascent proteins. Whole proteins then would be isolated and applied to the antibody matrix. Developing the resulting signal would result in a protein expression profile, which is useful in

essentially the same manner as the nucleic acid profile. A protein matrix could be used, for example, in evaluating antibody responses to pharmaceutical agents in order to eliminate possible cross-reactivity.

Moreover, where nucleic acids are used in the matrix, it is often beneficial to use variants (as defined below) of the molecules described hereinin. This can be used to account for genetic variations that are of little or no consequence to the function of the resultant gene product. Hence, they can account for wobble or conservative amino acid variations that do not perturb function, like variations in some of the protein motifs elucidated below. Thus, each position in the matrix can employ multiple nucleic acid probes that account for a series of variants.

Expression profiling may also be done in another embodiment using two-dimensional protein gels in which the inventive proteins are detected. The resultant profiles can be used in the same way as described.

Matrices useful for profiling may be constructed based on different criteria. Of course, the more relevant profiles will take into account expression of most human genes, preferably all of them. In certain situations, however, it is advantageous to look at a smaller subset. For example, if one were concerned about fetal neural toxicity, a fetal brain-specific matrix might be chosen. On the other hand, if one were interested in targeting mammary carcinoma tissue, a corresponding matrix could be used. Thus, matrices may be constructed using all of the sequences available from a tissue-specific library.

* * *

The following discussion relates to some of the various functional and structural groupings that would be of interest to the artisan wishing to construct profiling matrices. Of course, the artisan will also recognized that these functional descriptions may find additional applicability in the therapeutic and diagnostic applications discussed below.

Cell Cycle

15

20

25

A proliferating cell must coordinate replication and chromosomal separation to ensure that the genome is replicated

completely, and that a single copy is correctly inherited by each daughter cell. The cell cycle is the coordinated series of events that achieves these aims. Many of the key events are initiated by a family of conserved Seiren/threonine protein kinases, the cyclin-dependent kinases (CDKs), that are activated by the cyclin family of proteins (cyclins A-H). In turn, the cyclin-CDK complexes are modulated by other protein kinases or phosphatases, and by binding specific inhibitor proteins. The enormous variety of ways in which CDK activity can be regulated allows the cell to respond to internal signals generated by preceding events in the cell cycle and to external growth signals.

The somatic cell cycle is divided into four phases: DNA replication (S phase) and chromosome separation (M phase) are separated by gap phases (G1 and G2). At specific control points the decision to begin the next stage (DNA synthesis or mitosis) is carefully regulated.

15

20

25

30

35

Cdc2: the primary kinase: is especially required for the Gl-S transition and S phase. Cdc4 and Cdcb are involved at the restriction point: where the cell can decide to proliferate or arrest (Gl<->GD) and Cdc7 is a CDK activating kinase (CAK) as well as a subunit of TFIIH.

The Cyclin-CDK complexes are regulated in various ways. One is through phosphorylation by CDK activating kinases (CAK), like the YL5 kinase (Weel) and dephosphorylation by CDK associated phosphatases (CAP), like Cdc25A a member of the Cdc25 family (Cdc25A, B and C).

An other way of regulation occurs through two classes of CDK inhibitors (CKI), the INK4 proteins pl5, pl6, pl6, and pl9, who negatively regulates the cyclin D CDK complexes and second the p2l family with p2l, p27, and p57.

The cell cycle is also regulated through ubiquitin-mediated proteolysis involving the destruction of both cyclins and CDK inhibitors by the 26S proteasome, that requires an ubiquitin conjugating enzyme (UBC) and an ubiquitin ligase. The instability is conferred by PEST regions (cyclin D and E) or a ten amino acid

region in the amino terminus (degradation box) in the A- and B- type cyclins.

All these modifications play an important role for the cellular localization, because only the nuclear CDK-cyclin 5 complexes are functional for cell cycle. During GL phase of the cell cycle, cyclines A, E and D are synthesized and bind to their cyclin-dependent kinase (CDK) partners. CDK complexes containing cyclins A, E and D1 are then imported into and concentrated within nuclei. Cdkb- cyclin D3 has been localized to both 10 cytoplasmic and nuclear compartments, although only the nuclear complex is active. As cells enter S phase, cyclin A and cyclin E complexes remain within the nucleus, whereas cyclin DL relocalizes to the cytoplasm for proteolysis at the onset of S phase. Like Cdk2-cyclin A, Cdc2-cyclin A is nuclear and remains so until it is degraded during mitosis. By contrast, as a result of ongoing nuclear import and more rapid re-export, cyclin Bl, which binds to Cdc2 upon synthesis during S phase, is predominantly cytoplasmic. Cdc2-cyclin B2 is also cytoplasmic. although this might occur through anchoring of the complex to 20 some cytoplasmic constituent. At prophase, phosphorylation of cyclin Bl promotes accumulation of Cdc2-cyclin Bl in the nucleus. whereas cyclin B2 remains in the cytoplasm until nuclear envelope breakdown.

Two crucial regulators of Cdc2-cyclin B-Weel and Cdc25C

25 exist and are responsible for the G2 to M control point. Weel is
a nuclear protein throughout the cell cycle, whereas Cdc25C binds
to 14-3-3 proteins during interphase and remains predominantly
cytoplasmic. In some systems Cdc25C, like cyclin Bl, rushes
precipitously into the nucleus just before entry into mitosis.

The 110-kDa retinoblastoma (tumor suppressor) protein (RB), a pRB-family member is an important regulator of cell-cycle progression and differentiation. Like the E2F family (E2F1-5) or DP family (DP1-3) of transcription activators, RB suppresses inappropriate proliferation by arresting cells in G1 by repressing the transcription of genes required for the transition into S phase. Before the cell proceeds into S phase, RB becomes phosphorylated at multiple sites by the cyclin dependent protein

kinases (CDKs) and loses its transcriptional repressing activity. Phosphorylation of RB during late G1 phase results in the dissociation of the E2F-RB repressor complex which allows S-phase specific genes to be transcribed. Cyclin E is the evolutionary conserved target for E2F and interacts together with CDC2 in late GL-

For a proliferating cell it is vital that only undamaged DNA is replicated because if DNA damage is substantial, its replication can lead to chromosome loss or rearrangement. Thus, we find a G1<->S checkpoint in late G1 that requires tumor suppressor p53. A p53-dependent G1 arrest is effected by the cyclin dependent kinase inhibitor p21 through higher expression levels that inhibits almost all cyclin CDK complexes.

10

15

20

25

The kinase responsible for phosphorylating the unidentified kinetochore component in metaphase may be a member of the MAP kinase family and appears to be the proto oncogene c-MOS1 a cytostatic factor (CSF) in meiosis.

Several categories of proteins are coded for by clones of the invention within the overall group of "Cell cycle"and include, among others, the following:

PAZE-TZ protein: PAZE-TZ is a p53 responsive gene. The protein is predominantly expressed in brain, breast and kidney and represents a novel regulator of cellular growth. Isoforms are differentially induced by genotoxic stress (UV, gamma-irradiation and cytotoxic drugs)in a p53-dependent manner. The p53 tumor antigen is found in increased amounts in a wide variety of transformed cells. The protein is also detectable in many actively proliferating, nontransformed cells, but it is 30 undetectable or present at low levels in resting cells. P53 is postulated to bind as a tetramer to a p53-binding site (PBS) and to activate the expression of adjacent genes that inhibit growth and/or invasion. Deletion or inactivation of one or both p53 alleles reduces the expression of tetramers, resulting in decreased expression of the growth inhibitory genes. This mechanism is found in tumors of several types. (OMIN *191170) Clones in this category include: amy2_121m2

Cell structure and motility

One of the major differences between prokaryotes and eukaryotes is the ability of the eukaryotic cell to adopt very different shapes dependent on its function during the differentiation process. Animal cells vary from being round to extended cylindric forms like motorneurons or muscle cells. In humans, more than 100 different cell types can be distinguished, each having a characteristic shape. The form of a cell often is 10 closely related to its capacity to move. Some completely differentiated cells like fibroblasts can still change their form actively, thereby migrating. Other cell types serve as motor elements - "macroscopically" like muscle cells or "microscopically" like ciliated epithelia. Such tasks are fulfilled by a big class of proteins; on the one hand responsible - 15 for maintenance of cell structure and contacting neighbor cells or the intercellular matrix and on the other hand for cell motility. These topics cannot be regarded separately: The motility apparatus e.g. must be fixed in the cytoskeleton. Three: 20 different types of filaments can be distinguished: Actin filaments, tubulin filaments and intermediate filaments, each present in almost all types of cells-

Actin filaments (F-actin) are built up of monomers (G-Actin). In muscle cells, actin, myosin, for both of which several paralogous genes are known, as well as many more proteins are constituents of the contractile apparatus.

The "thin" and "thick filaments" in a muscle cell consist mainly of actin and myosin, respectively.

Several different proteins are responsible for the anchoring 30 of the actin filaments in the Z-disks (e.g. alpha-actinin and desmin) or at the end of the myofibers in the cell membrane.

Troponin I, $-C_1$ -T and Tropomyosin - associated with actin - confer the Ca++- dependent triggering of contraction.

Length of the sarcomere is controlled by the giant protein titin.

In smooth muscle, there is no troponin. Contraction activity is controlled by phosphorylation / dephosphorylation of myosin by a specialized kinase instead. Contractile fibers are not organized in sarcomeres.

Apart from contributing to muscle contraction, the actomyosin system is responsible for many other motions at cellular level, e.g. the amoeboid movement of pseudopodia or the fission of cells at the end of mitosis by a contractile ring.

10

15

20

25

30

35

Besides this, actin fibers fulfill structural tasks like maintenance of the shape of stereocilia or microvilli. Here, actin filaments are connected by proteins like fimbrin. But not only specialized structures like the mentioned ones contain actin fibers. There is a network covering the complete cell volume with F-actin as a major constituent. Whereas the actin filaments in the structures mentioned above are relatively stable, this F-actin is highly dynamic. Management of the network structure and turnover is achieved by connecting proteins like alpha-actining fimbrin or fill-in; turnover is regulated by gelsoling villing and different capping—and fragmentation—proteins.

Microtubules are built up of alpha-beta tubulin heterodimers. Turnover of filaments is achieved by building-in and releasing of monomers with different time constant rates at both ends. The resulting cycle is called "treadmilling". Thirteen strings of tubulin duplets build up one subfiber, whereas one fiber contains two or three of those. A complete axoneme consists of 9 radial and 2 central fibers. This "9+2" - structure is the basis both of flagella, their basal bodies and centrioles. In flagella, several additional structures like radial elements exist. Nexin connects the fibers and dyneine is the motor ATPase which shifts the fibers relative to each other. Several genetic diseases like the Cartageneric syndrome are caused by deficiencies of distinct proteins in cilia.

Besides this, microtubules are abundant in all types of cells. They are part of a delivery system for organelles, e.g. in

WO 01/98454

15

25

30

35

PCT/IB01/02050

the golgi apparatus. A further very important system based on microtubules is the mitotic spindle, it is organized by the centrosomes. Besides many other components, the major part of a centrosome are two centrioles which are built up of nine microtubule-triplets. Most remarkably, new centrioles are not synthesized de novo but generated by duplication of old ones.

Cytoplasmic microtubules are associated with many different. proteins. Two major classes are known: The MAPs ("microtubuleassociated proteins", with molecular masses between 200 and 300 kD) and the much smaller tau-Proteins with a MW between 60 and 70 kD. These proteins regulate the treadmill-process and the interaction with other structures in the cell-

Besides actin and myosin the so-called intermediate filaments constitute a third class of filaments. In contrast to the former two groups, they do not participate in motility, nor are they dynamic structures subject to a vivid turnover. The most important ones are neurofilaments (in neurons), keratin filaments (mainly in epithelial cells), and vimentin filaments (in many sorts different cell types).

The biological function of both the cytoskeleton as well as 20 contractile apparatus of a cell does not end at the cell membrane. Cells must be embedded in the extracellular matrix, all cells of a muscle must act as one single mechanical unit and epithelia must resist macroscopic mechanical forces. Hence, cell adhesion and the extracellular matrix are closely connected to the cytoskeleton. Vincullin is one of the proteins which serve as an anchor for intracellular fibers (actin). Different types of desmosomes and tight junctions connect neighbor cells with intercellular fibers. On the inside, cytoplasmic plaques connect them to the cytoskeleton. These structures, on the one hand, serve as mechanical elements whereas gap junctions, on the other hand, connect cells metabolically.

The extracellular matrix consists of a network of proteins, glycoproteins and polysaccharides. Different proteins are present in relation to different mechanical demands:- Elastin is found in tissues with high elasticity (lungs, heart) whereas collagen,

a more hard-wearing protein, is found in tendons and ligaments. Fibronectin is an extracellular protein highly important for cell adhesion.

Reference: Murray J et al (1992): Cell Motil Cytoskeleton 5 22: 211-223.

Within the overall group of Cell Structure and Motility several categories of proteins are coded for by clones of the invention:

Ankyrins: Ankyrins are peripheral membrane proteins which interconnect integral proteins with the spectrin-based membrane skeleton. Thus these proteins are involved in coupling of cyto skeleton and cell membrane. OMIN reports that Ankyrins have associations (as potentially diagnostic, therapeutic, causative, and/or related, etc...) with the following diseases: 1) Heriditary Spherocytosis (OMIN *182900); 2) Hemolytic Poikilocytic Anemia due to reduced ankyrin binding sites (OMIN 141700); 3) Atypical Elliptocytosis (OMIN 225450); 4) Autosomal recessive spherocystosis (OMIN #270970); 5) Werner Syndrome (OMIN *277700); and b) Rhesus-unlinked type Elliptocytosis (OMIN #130600). Ankyrin bindung glycoprotein proteins mediate Ankyrin effects, especially in neuronal adhesion and prostate tumour vcell transformation: Clones in this category include: amy2_121f19.

Tropomyosins are ubiquitous proteins of 35 to 45 k) associated with the actin filaments of myofibrils and stress fibers. They are involved in cardiomyopathies (OMIN *191030, *191010, *190990, *600317). Clones in this category include: tes3_165.

Differentiation/Development

10

15

20

25

Almost every multicellular organism originates from meiotic cell divisions and the recombination of a paternal and a maternal set of chromosomes. After fertilization of the egg, all cells of a body originate from this one cell. Thus the cells of the developing body are initially genetically alike. But

35 phenotypically they become very different. They are specialized to a certain cell type and arranged in an organized pattern to a certain type of tissue and the whole structure has the well-

WO 01/98454

PCT/IB01/02050

defined shape of an organ. All these features are determined by the DNA sequence of the genome, which is reproduced in every cell. Each cell acts on the genetic instructions given to a certain time and at a certain place of development and plays its individual part in the multicellular organism. Cell differentiation may be divided into three general steps: cell cycle exit, apoptosis protection and tissue specific gene expression. These processes are coordinated to provide the final and unique tissue characteristics.

10 An animal cell that has achieved a certain level of development is said to be determined. This differentiation of a cell may be irreversible and in that case the cell may be renewed only by simple duplication. Other cells are renewed by means of stem cells which are immortal (e.g. stem cells of the bone 15 marrow, epidermal stem cells). The genetic control of development is extensively studied in non-vertebrates and vertebrates. The classical animal model is the fruit fly Drosophilia and the modern model is the transgenic mouse. Animal transgenesis has proven to be useful for physiological as well as 20 physiopathological studies. Besides the approach based on the random integration of a DNA construct in the mouse genome, gene targeting can be achieved using totipotent embryonic stem cells for targeted transgenesis. Transgenic mice are than derived from the embryonic stem cells. This allows the introduction of null mutations in the genome (so-called knock-out) or the control of the transgene expression by the endogeneous regulatory sequence of the gene of interest (so-called knock-in). Mice can be created that express wild-type genes, mutant genes, marker genes or cell lethal genes in a tissue specific manner. These animal 30 models allow to follow changes in tissue and organ development and lead to a better understanding of the cellular function of many genes or to the generation of animal models for human diseases. Fundamental problems in immunology, onset and development of cancer, regulation in fatty acid metabolism, aspects of cardiovascular function, control of the central 35 nervous system development, analysis of reproductive development and function are only some examples of research interests.

The final stage of cell differentiation is growth arrest. In animal tissues with rapid cell turnover terminally differentiated cells undergo programmed cell death. The cells have the ability to kill themselves by activating an intrinsic cell suicide program when they are no longer needed or have become seriously damaged. The execution of this program is termed apoptosis. Apoptosis is of importance for development and homeostasis of animals. The key components of this program have been conserved in evolution from worms (C. elegans) to insects (Drosophilia) to humans. The roles of apoptosis include the 10 sculpting of structures during development, deletion of unneeded cells and tissues, regulation of growth and cell number, and the elimination of abnormal and potentially dangerous cells. In this way apoptosis provides "quality control mechanism" that limits the accumulation of harmful cells, such as virus-infected cells and tumor cells. On the other hand inappropriate apoptosis is associated with a wide variety of diseases, including AIDS, neuro-degenerative disorders and ischemic stroke. Because it is now clear that apoptosis is a result of an active, gene-directed 20 process, it should be eventually possible to manipulate this form of cell death by developing drugs that interact with its recently identified mechanisms of action. Inducers of cell differentiation, cell cycle arrest and apoptosis might be the novel molecular targets for new anticancer agents in addition to the signaling pathways for growth factors and cytokines. 25

<u>Proteins, factors, receptors and genes of importance in apoptosis:</u>

Proteases:

35

- Calpain, an intracellular cysteine protease, exact role 30 unknown.
 - Caspase-1 to Caspase-11, a family of proteases synthesized as an inactive proenzyme. Targets of the activated enzymes include: poly(ADP-ribose) polymerase, DNA-dependent protein kinase, U1 ribonucleoprotein, nuclear laminins and cytoskeleton components (actin).

- Granzyme B, a serine protease released by cytotoxic T-

Receptors:

- CD 95 (synonyms: Fas, APO-1), a receptor protein of the 5 TNF-receptor family which includes TNF-R1 and TNF-R2 with the common characteristic of a 7D amino acid cytoplasmic domain.
 - FADD (synonym: MORT-1), a cytoplasmic protein
 - DR-3 (synonym: APO-3) a member of the TNF-receptor-family
 - DR-4 and DR-5

10 Genes:

- ced-3, ced-4 and ced-9 encode the general apoptotic and antiapoptotic program in Caenorhabditis elegans. Apaf-3 is the mammalian homologue of ced-3.
- Bcl-2 / Bcl-xL / Bax / Bcl-xS / Bak: a large gene family that can either inhibit or promote apoptosis.
 - Cytokine response modifier A₁ a cowpox virus gene whose gene product inhibits caspases.

Others:

- Caspase-activated DNase (CAD) and its inhibitor (ICAD). 20 causes DNA fragmentation in the nucleus
 - Ceramide, a complex lipid that acts as a second messenger.
 - c-Jun N-terminal kinase (JNK) is a proline-directed kinase
 - p53 protein, is essential for the induction of apoptosis as a response to chromosomal damage.
- 25 RAIDD, a death signal-transducing protein.
 - Receptor interacting protein (RIP) is an accessory protein with a death domain and a serine/threonine kinase activity.

- Sphingomyelinase, an enzyme that hydrolyzes the complex lipid sphingomyelin to ceramide.

- Tumor necrosis factor (TNF) is a type -II membrane protein
- TNF-receptor associated factor (TRAF2), is an accessory protein that can bind to both TNF-R1 and TNF-R2.

Within the overall group of Differentiation/Development, several categories of proteins are coded for by clones of the invention:

Notch family proteins: Notch family molecules are negative regulators of neuronal differentiation in early brain development. Clones in this category include: amy2_li24.

Testis-specific Y-encoded proteins: The TSPY genes are arranged in clusters on the Y chromosome of many mammalian species. TSPY is believed to function in early spermatogenesis and is a candidate for GBY, the putative gonadoblastoma-inducing gene on the Y. These proteins are involved in early spermatogenesis. Clones in this category include: amy2_7j5.

Inflammation-mediating proteins: Inflammation is a basic mechanism responsible for recruiting and activation of immuncompetent cells. By various mediators, cells are activated and triggered to differentiate. Hyperactivation of these pathways leads to various disease states: In neuronal tissues, in inflammatory diseases such as experimental autoimmune encephalomyelitis (EAE), neuritis(EAN) and uveitis (EAU) allograft inflammatory factor-1 is produced by macrophages and microglia cells. Clones in this category include: amy2_2b19.

Intracellular transport and trafficking

15

30

35

Eukaryotic cells rely for their viability on the partitioning of many basic cellular processes into membrane-bounded organelles. These are the nucleus, endoplasmic reticulum (ER), Golgi apparatus, endosomes, lysosomal compartments, mitochondria and peroxisomes. Most molecules destined for the lysosome, cell surface and outside the cell are routed through

the ER and Golgi, which together with the vesicular intermediates between them, comprise the secretory pathway (Palade 1975). In the ER and Golgi compartments proteins are sorted, modified and often assembled into complexes en route to their final destination. Incorrectly assembled proteins are retained in the ER until they fold correctly or are targeted for degradation. Additional proteins are translocated into and function within the lumenal spaces of organelles or are secreted. Thus a large proportion of proteins synthesized require targeting to membranes 10 either for insertion into or transport across them. A major purpose of this is growth. The secretory pathway is dependent on an intact cytoskeleton and also closely linked to general metabolism by affecting ribosome biogenesis (Mizuta and Warner, 1994). A huge number of proteins is required for targeting. 15 translocation and sorting of newly synthesized proteins.

The first step in sorting is the recognition of cis-acting targeting or signal sequences that organelle-targeted proteins contain. This is carried out by cytosolic targeting factors and/or receptors on the membrane to which the protein is targeted. In some cases the primary sequences are extremely degenerate, with only the overall character being conserved (hydrophobicity for an ER signal sequence, helical amphiphilicity for mitochondrial targeting sequence (Kaiser et al., 1987; Lemire et al., 1989). Following the targeting step, proteins are either inserted into or transported across the membrane (translocated) through a proteinaceous apparatus (termed the translocon). The translocon include or recruit motors to drive the translocation process in the correct direction (Schatz and Dobberstein, 1996).

Defined intracellular protein transport steps:

30 • ER

20

25

35

- targeting to the ER
- translocation into the lumen of the ER, and, depending on the presence of certain signals in the peptide sequence transport through the golgi complex
 - Mitochondria
 - targeting
 - translocation
 - Peroxisomes

- The general secretory pathway
- protein modification, assembly and quality control in the ER
 - vesicle-mediated trafficking
- vesicle docking and fusion
 - transport through the golgi apparatus and sorting at the trans-golgi
 - transport to the cell surface
 - transport routes to the lysosome
- 10 Endocytosis

5

20

25

35

- Specialized protein transport routes
- Protein export from the cytoplasm

References: Palade, 6 (1975) Science 189:347-358; Mizuta et al. (1994) Mol Cell Biol 14: 2493-2502; Kaiser et al. (1987) 15 Science 235: 312-317; Lemire et al. (1989) J Biol Chem 264: 20206-20215; Schatz et al. (1996) Science 271: 1519-1526.

Rab proteins

In eukaryotic cells the compartmentalisation of processes is prerequisite for a tight regulation of processes activities. The cells contain a highly dynamic set of membrane compartments that are responsible for packaging, and recycling proteins and other molecules. Trafficking between organelles within the secretory occurs as vesicles derived from a donor compartment fuse with specific acceptor membranes, resulting in the directional transfer of cargo molecules. This process is tightly controlled by the Rab/Ypt family of proteins (reviewed by Novick and Zerial, 1997), a branch of the superfamily of small GTPases. proteins regulate a variety of functions, including vesicle 30 translocation and docking at specific fusion sites. Rabs may also play critical roles in higher order processes such as modulating the levels of neurotransmitter release in neurons, a likely mechanism in synaptic plasticity that underlies learning and memory (Geppert and Sudhof, 1998)-

Small GTPases share a common three-dimensional fold that, in the GTP bound state, can bind a variety of downstream effector proteins. GTP hydrolysis leads to a conformational change in the "switch" regions that renders the GTPase unrecognizable to its

effectors. In this way, by localizing and activating a select set of effectors, a common structural motif is used to control a wide array of distinct cellular processes.

The final steps in membrane fusion are likely to be driven by a set of proteins known as SNAREs. After a vesicle becomes cytoplasmic domains of VAMP (also synaptobrevin) and syntaxin on opposing membranes, in combination with a SNAP-25 molecule, coalesce into an elongated -helical bundle (Poirier et al., 1998; Sutton et al., 1998), which may lead to fusion. Because numerous SNARE isoforms have been identified that localize to distinct membrane compartments, it was originally proposed that the specificity of interaction between the SNARE proteins accounted for the specificity in membrane trafficking. Recent results, however, suggest that SNAREs are not specific in their ability to form complexes in suggesting that trafficking specificity additional factors (Yang et al., 1999). In this regard, Rab proteins are strong candidates for governing the specificity of vesicle trafficking. Like the SNAREs, many isoforms (40) of the Rab family have been identified that localize to specific. membrane compartments (reviewed by Novick and Zerial, 1997).

10

15

20

25

30

35

Concomitant with the SNARE cycle, Rab proteins undergo a intricate cycle of membrane and protein interactions. Rabs are posttranslationally modified at C-terminal cysteines by the addition of two geranylgeranyl groups, which mediate membrane association when the Rab is in the GTP-bound state. After guanine nucleotide hydrolysis occurs, the Rab is extracted from the membrane upon forming a complex with a cytosolic GDP-dissociation inhibitor (GDI). This cytosolic intermediate is then recycled onto a newly forming vesicle, most likely through a secondary factor termed a GDI dissociation factor (GDF), which displaces GDI. After the Rab becomes membrane bound, a guanidine nucleotide exchange factor (GEF) promotes release of GDP and the subsequent loading of GTP. In its GTP-bound conformation, the Rab is then free to associate with its specific set of effectors; which can in turn trigger events leading to the eventual fusion of the vesicle with a target membrane. To complete the cycle, perhaps after or concurrent with membrane fusion, a GTPase activating protein (GAP) accelerates nucleotide hydrolysis, switching off

the GTPase. The remaining GDP-bound Rab can then participate in a new round of fusion.

Rab interactions with effectors are likely to regulate vesicle targeting and membrane fusion in three ways. First, a Rab may specifically facilitate vectorial vesicle transport. Vesicles transported from their site of origin to compartments likely through associations with cytoskeletal elements and transport motors. A protein has been identified with , a domain structure that suggests a connection between the cytoskeleton and the Rabs. This protein, called Rabkinesin-L, contains a kinesin-like ATPase motor domain followed by a coiledcoil stalk region and a RBD that specifically binds Rabb (Echard et al., 1998). An additional link with the cytoskeleton is provided by the Rab effector, Rabphilin-3A. Rabphilin-3A has been shown in vitro to interact with -actinin, an actin-bundling protein, but only when not bound to Rab3A (Kato et al., 1996). These results raise the intriguing possibility that Rab proteins regulate vesicle interactions with the cytoskeleton and thereby play an active role in targeting vesicles to their appropriate destinations.

10

15

20

25

35

Second, Rab proteins may regulate membrane trafficking at the vesicle docking step. A number of Rab effectors, including Rabaptin-5, EEAL, Rabphilin-3A, and Rim, may serve as molecular tethers. Each effector protein contains a RBD, followed by a linker region (some having the potential to form elongated coiled-coil structures), and a domain capable of interacting with a second Rab or the target membrane. Rabaptin-5, for example, contains two RBDs, one near the N terminus that specifically recognizes Rab4 and a second near the C terminus that binds Rab5 30 (Vitale et al., 1998). Both Rim, which is localized to the target membrane, and Rabphilin-3A, which is localized to the vesicle, contain N-terminal RBDs and C-terminal Ca2+-binding C2 domains. implicating these effectors in synaptic localization or docking in response to Ca2+ influx (Wang et al., 1997). Tethering effectors may also recognize protein complexes on the acceptor membrane. Sec4p, a yeast Rab3A homolog, interacts with the exocyst (Guo et al., 1999), a complex of seven or more subunits that is assembled at sites of vesicle fusion along the

plasma membrane. The exocyst complex may therefore function as a landmark for Rab/effector-mediated vesicle docking.

Third, once a vesicle has become tethered to its fusion site, Rab proteins may selectively activate the SNARE fusion machinery. The mechanism of this activation is unknown but may involve direct interactions of Rabs or, more likely, their effectors with SNAREs. For example, Hrs-2 is a protein that binds to SNAP-25 and contains a Zn2+-finger motif characteristic of Rab-binding proteins such as Rabphilin-3A, Rim, EEAl, and Noc2, suggesting that Hrs-2 may form a physical link between Rabs and SNAREs (Bean et al., 1997). In addition, certain mutations in the syntaxin-binding protein Slylp, the Seclp homolog utilized in ER to Golgi trafficking, eliminate the requirement for Yptlp, a Rab protein that functions at this trafficking step (Dascher et al., 1991). Rabs may therefore regulate SNARE associations through Secl family members. In support of this idea, a Rab effector was recently found to interact with a vacuole Rab, a Seclp homolog, and a SNARE protein (Peterson et al., 1999), which suggests that this effector serves to connect Rab and SNARE function. In this way, Rabs and their effectors may facilitate the correct pairing 20 of SNAREs.

10

25

30

35

References: Dascher et al. (1991) Mol. Cell. Biol. 11, 872-885; Echard et al. (1998). Science. 279; 580-585; Geppert et al. (1998) Annu. Rev. Neurosci. 21, 75-95; Guo et al. (1999). EMBO J. 18, 1071-1080; Kato et al. (1996) J. Biol. Chem. 271, 31775-31778; Novick et al. (1997) Curr. Opin. Cell Biol. 9, 496-504; Peterson (1999) Curr. Biol. 9, 159-162; Poirier et al. (1998) Nat. Struct. Biol. 5, 765-769; Vitale et al. (1998) EMBO J. 17, 1941-1951; Wang et al. (1997) Nature. 388, 593-598; Yang et al. (1999) J. Biol. Chem. 274, 5649-5653.

Within the overall group of Intracellular Transport and Trafficking several categories of proteins are coded for by clones of the invention.

Vesicular trafficing: Various proteins are involved in trafficing of vesicles inside the cell and for the exocytotic pathway. For example, Sec7 of Saccharomyces cerevisiae takes function in vesicular traficking. Synaptotagmins are essential for Ca(2+)-regulated exocytosis of neurosecretory vesicles. Other proteins such as Dynamin are microtubule-associated force-

producing proteins, which are involved in the production of microtubule bundles. By binding and subsequent hydrolysation of GTP such proteins provide the motor for vesicular transport during endocytosis. Clones in this category include: amy2_14b5, amy_2013 and fkd2_3k1.

<u>Protein sorting:</u> Protein sorting is a process essential for the maintenance of a cells functionality and structural integrity. Most proteins perform their biological function in special compartments in the cell. The process of sorting is complex and highly regulated. Clones in this category include: mel2_7g14.

Metabolism

10

15

20

25

30

35

This group includes proteins which are involved in the uptake and consumption of nutrients, and enzymes which are part of the biochemical pathways for energy metabolism or which are involved in the supply of building blocks of nucleic acids, proteins (NTPs, dNTPs, amino acids) for DNA/RNA and protein synthesis, and fatty acids (membranes), to allow for the generation of higher order structures. This group constitutes the most important and largest group in prokaryotes and lower eukaryotes. The higher the evolutionary level of an organism is, however, the more other protein classes like 'signal transduction', 'cell cycle' and 'differentiation and development' increase in importance and number of representatives.

Proteins involved in the metabolism of energy and compounds (here: other than nucleic acids or proteins) are usually the products of house keeping genes, they are often constitutively and/or ubiquitously expressed.

Several categories of proteins are coded for by clones of the invention within the overall group of Metabolism:

Fatty acid metabolism: OMIN lists more than 50 diseases caused by pathologic altered fatty acid metabolism. L-acyl-glycerol-3-phosphate acyltransferase is involved in fatty acid metabolism and is ubiqitous expressed, with a slight predominance in uterus, placenta and foreskin. Clones in this category include: amy2_2c22

Repair and surveillance of protein damage: Several classes of protein are involved in reapair and surveilance of protein damage. L-isoaspartyl methyltransferase (Pimt), as an example, is a highly conserved enzyme utilising S-adenosylmethionine (AdoMet) to methylate aspartate residues of proteins damaged by agerelated isomerisation and deamidation. Clones in this categroy include: fbr2_78i21.

Nucleic acid management

10

15

30

35

:1

The genetic information is stored in the form of nucleic acids in all organisms. Two kinds of nucleic acids exist, DNA and RNA. Whereas the more stable DNA in most organisms constitutes the storage form of the genetic information, the labile RNA and in particular mRNA is an intermediate used for the temporal expression of specific genes.

In eukaryotes, DNA is usually a double stranded linear molecule consisting of two antiparallel strands and made up of a deoxyribose, a phosphorus backbone and the four bases A, C, G, and T. The DNA of some organisms has a ring structure. The structure of DNA was unraveled years ago by Watson and Crick. DNA is directional molecule determined by the C-atoms of the sugar.

The most important processes dealing with nucleic acids are:

- replication (e.g. DNA polymerases, Telomerase)
- transcription (RNA polymerases)
- RNA processing (maturation splicing and degradation)
- in addition, enzymes and proteins exist which require a nucleic acid (mostly RNA) in the active center to be functional (ribozymes e.g. RNase, Ribosomal proteins)

The DNA of a cell is replicated in the S-phase of the cell cycle. Several enzymes carry out the task of doubling this nucleic acid. As all steps of the cell cycle, also the process of replication is tightly regulated. The enzyme DNA polymerase and several other proteins are involved in this process. Whereas many prokaryotes do have only one origin of replication (i.e., the starting point of the replication cycle), in eukaryotic DNAs (chromosomes) multiple such start points exist. The switch from the synthesis (S) phase to the subsequent G2 or M phases of the cell cycle are dependent on the completion of the replication.

This makes clear, that a number of proteins are involved in the replication itself as well as in the control of the process. Since most eukaryotic chromosomes are linear structures, additional proteins and enzymes are necessary to make sure that the structure is maintained through successive generations. This includes those proteins necessary to build the three dimensional structure of chromosomes (e.g. histones) and the structural network of the nucleus and nucleolus (including the defined localization of transcriptionally active genes in the vicinity of nucleoli) but also such enzymes as telomerase which guarantees the integrity of the chromosomal ends.

10

15

The expression of genes is usually performed in two steps. First a messenger RNA (mRNA) is produced (transcribed) in one to many copies and second this mRNA is translated into the protein product. The regulation of transcription is discussed under the separate heading 'transcription factors', but also the classes 'signal transduction', 'development', 'cell cycle' and others are affected as the expression of certain genes determines the fate of a cell or organism.

20 The primary transcript (hnRNA - heterogeneous nuclear RNA) is a single stranded one-to-one copy of the gene as it is located on the chromosome. Before a protein can be translated, already during transcription the process of maturation is initiated. Firstly, a 5' cap structure is enzymatically and covalently added 25 to the RNA, blocking the 5' end of the RNA. Second, when the RNA polymerase has terminated polymerization, the enzyme poly A polymerase adds varying numbers of adenine residues to the 3° end of the transcript. This enzyme recognizes the sequence AAUAAA or AUUAAA (+ some minor variations), cuts the RNA 10 - 30 30 nucleotides downstream and adds the A residues. The size of the poly A sequence affects the stability of the RNA. Finally, in the process of splicing, the introns present on the genomic level and also present in the hnRNA are spliced out by a multi-protein complex consisting of several proteins and RNAs. The finally maturated mRNA is exported to the cytoplasm where it is 35 translated with help of the ribozymes.

The half life of RNA is usually much shorter than that of DNA. Usually, the mRNA is degraded shortly after synthesis, to guarantee a very defined window of expression of a given gene.

This regulation is necessary to specifically maintain or change the set of proteins present at any time in a cell. Specific regions in the 3'UTR (untranslated region) determine the stability of the mRNA in the cytoplasm before it is degraded by RNases, enzymes consisting both of protein and RNA.

References: Watson and Crick (1953) Nature 171: 737-738.

Several categories of proteins are coded for by clones of the invention within the overall group of "Nucleic acid management"and include, among others, the following:

10

15

20

25

30

35

Proteins induced by DNA-Damage: There are several distinct pathways responsible for repair of DNA. Nucleotide excision repair is the most versatile DNA repair pathway and isthe main defense of mammalian cells against UV-induced DNA damage. Defects in proteins involved in this pathway can lead to inherited disorders (such as xeroderma pigmentosum OMIN *278700; *278720; *278740 and *194400; Cockayne's syndrome OMIN *216400 and trichothiodystrophy OMIN #601675). Study of UV-sensitive yeast RAD mutants has greatly aided this process and has revealed strong conservation of the components of nucleotide excision repair in eukaryotes. Clones in this category include: amy2_11n4 and tes3_10i16.

Proteins involved in Loading of transferRNAs: transfer RNAs must be coupled to an aminoacid, which then is transported to the peptideyl-transferase centre of the ribosome. Clones in this category include: fbr2_78cl2.

Cytosolic ribosomal proteins: Several proteins are part of the eukaryotic ribosomal peptidyl transferase center or modulate the activity of this centre. Such proteins can find application in modulation of ribosome assembly, maintenance and activity. Clones in this category include: amy2lil

<u>Histones:</u> Histones are DNA-binding protein responsible not only for DNA structure and folding and packing, but also are discussed to be involved in activation and silencing of large chromosomal regions. Clones in this category include: tes3_3lalD.

mRNA-binding proteins: mRNA-binding are involved in regulation of mRNA folding, translation and stability. For example, the VILIP protein binds specifically to the

3'untranslated region of the neurotropin receptor mRNA. Clones in this group include amy2_2gl2.

Signal transduction

Cells in higher order organisms need to continuously communicate with its environment especially with other cells of the same organism in order to maintain the function and specialization of the whole system these cells are part of. This important task of communication is performed with help of cellsurface receptors which receive and transmit signals from outside into the cell.

<u>G-proteins</u>

5

10

15

20

25

30

35

The largest known family of cell-surface receptors is that of the G-protein-coupled receptors, which mediate the transmission of diverse stimuli such as neurotransmitters, glycopeptides, hormones, peptides, odorant molecules, and photons. The functional unit of these receptors is composed of the receptor molecule itself (GPCR) which is anchored in the cytoplasma membrane with seven membrane spanning domains, the heterotrimeric G-protein which is composed of and -subunits (G and G) and the effectors that interact with G and / or G . In particular, the dissociated G and G can regulate the activities of a number of effector molecules such as adenylate cyclases, phopholipase C isoforms, ion channels, and tyrosine kinases, resulting in a variety of cellular functions. The process of signal transduction must be tightly regulated and reversible in order to avoid overstimulation, to achieve signal termination, and render the receptor responsive to subsequent stimuli [[acovelly L. et al., (1999) FASEB J. 13, 1-8, Hamm, H.E. (1998) J. Biol. Chem. 273, 669-6721.

G-proteins are GTPases that, upon binding of GTP change their conformation which in return unmasks structural motives, in particular the so called effector loop, which can mediate the interactions to target proteins, or effectors, for the GTPases. This ability enables the GTPases to cycle between active, GTP-bound and inactive, GDP bound conformations and in the process to function as molecular traffic lights in a multitude of signal transduction pathways. The most important of these signal transduction pathways that are regulated with help of G-proteins

are that of the phospholipase (/ protein kinase (and that of the adenylate cyclase / protein kinase A.

The cycling of GTPases is tightly regulated by three main classes of proteins: The exchange of hydrolyzed GDP for a fresh GTP is facilitated by guanosine nucleotide exchange factors (GEFs), the hydrolysis of GTP to GDP is sped up by GTPase-activating proteins (GAPs), and the dissociation of GDP from the GTPases is inhibited by GDP dissociation inhibitors (GDIs) ETapon and Hall (1997) Curr.Opin. Cell. Biol. 9, 86-92, Van Aelst and D-Souza-Schorey (1997) Genes Dev. 11, 2295-23221.

SOC-family

10

15

20

25

30

35

A conserved motif that was originally identified in proteins that negatively regulate the signaling action of cytokines was termed SOCS box, the Suppressor Of Cytokine Signaling. Based on homology, five distinct structural protein classes have been identified since that carry this motif. The function of most of these proteins is presently not known. Common to the proteins is only the SOCS box which is located near the C-terminus of the respective peptides. Recently, the SOCS box has been demonstrated to induce binding of proteins to elongins B and C which could target the proteins (and bound substrates) to the proteasomal protein degradation pathway (Kamura, T. et al. (1998) Genes Dev. 12, 3872-3881; Zhang, J.-G. et al. (1999) Proc. Natl. Acad. Sci. USA 96, 2071-2076).

The class where the SOCS box was originally described contains several members (SOCS-1-SOCS-7 and CIS). In addition to the SOCS box, these proteins also contain a SH2 (Src-homology 2) domain and a variable N-terminus. These SOCS proteins appear to form part of a classical negative feedback loop that regulates cytokine signal transduction. Upon cytokine stimulation, expression of SOCS proteins is rapidly induced and the proteins inhibit further cytokine action. The mode of action of the SOCS proteins is variable. While SOCS-1 binds and inhibits the JAK (Janus kinases) family of cytoplasmic protein kinases ENarahzaki M. et al. (1998) Proc. Natl. Acad. Sci. USA 95, 13130-13134. Nicholson, S.E. et al. (1999) EMBO. J. 18, 375-3851. CIS appears to act by competing with signaling molecules such as the STATs (Transducers and Activators of Transcription) family for binding

to phosphorylated receptor cytoplasmic domains [Yoshimura, A. et al. (1995) EMBO J. 14, 2816-2826; Matsumoto, A. et al. (1997) Blood 89, 3148-31541.

A second class of SOCS box protein contains additionally WD-4D repeats which were initially identified in the mouse WSB-1 and -2 proteins. The functions of WD-4D proteins are not completely understood but seem to be rather divergent. In Cdc4p the WD-4D repeats probably are necessary for binding the substrate for Cdc34p [Mathias N. et al. (1999) Mol. Cell Biol. 19, 1759-1767]. Cdc4p is a component of a ubiquitin ligase that tethers the

Cdc4p is a component of a ubiquitin ligase that tethers the ubiquitin-conjugating enzyme Cdc34p to its substrates. The posttranslational modification of a protein by ubiquitin usually results in rapid degradation of the ubiquitinated protein by the proteasome. The transfer of ubiquitin to substrate is a multistep process where WD-4O repeats might play an important function.

10

15

20

25

30

35

Other WD-4D containing proteins (e.g. the retino blastoma binding protein RbAp4B) have been shown to bind metal ions (Zinc) and that this metal binding might mediate and/or regulate protein-protein interactions which are functionally important in chromatin metabolism EKenzior, A.L. and Folk, W.R. (1998) FEBS Lett. 440, 425-4293. These proteins are involved in the RAS-cAMP pathway that regulates cellular growth EAch R.A. et al. (1997) Plant Cell 9, 1595-16063.

The SPRY domain has been identified in pyrin or marenostrinal a protein which is mutated in patients with Mediterranean fever and which is similar to the butyrophilin family. While butyrophilins seem to be involved in the lactation process in mammals, the function pyrin is unknown. Three proteins (SSB-L to -3) have been identified to contain both SPRY and SOCS box motifs. The function of these proteins is also not known.

Ankyrin repeat containing proteins share a 33-residue repeating motif, an L-shaped structure with protruding -hairpin tips which mediate specific macromolecular interactions with cytoskeletal, membrane, and regulatory proteins. These proteins play fundamental roles in diverse biological activities including growth and development, intracellular protein trafficking, the establishment and maintenance of cellular polarity, cell adhesion signal transduction, and mRNA transcription. Three proteins that

contain ankyrin repeats (ASB-1 to -3) have been identified to contain a C-terminal SOCS box additionally to the ankyrin repeats. The function of these proteins or the individual domains remains to be discovered [Hilton, D.J. et al. (1998) Proc. Natl. Acad. Sci. USA 95, 114-1191.

A few small GTPases (RAR and RAR like) do also contain a SOCS box. GTPases are involved in signal transduction during cellular communication. The function of the SOCS box in this type of proteins is currently unclear EHilton, D.J. et al. (1998) Proc. Natl. Acad. Sci. USA 95, 114-1191.

Ca 2+ as second messenger

The bivalent cation Ca2+ is, besides cAMP, one of the two major second messengers in eukaryotic cells. Its intracellular concentration is tightly regulated and usually kept very low compared to the cell's environment. Ca2+ binding proteins and transporters (Gap junction, Voltage-gated, second messengergated) help to sequester huge amounts of the ion in various organelles from where Ca2+ can be released upon extracellular stimuli. E.g. the contraction of the muscle is dependent on the presence of Ca²⁺ ions which are readily transported back into the organelles in order for the muscle to relax. In transduction, Ca²⁺ functions as a second messenger that activates Ca2+ dependent processes through the activation of Ca2+/calmodulin dependent protein kinases (CaM kinases) which are the major effector molecules of Ca2+. In the signaling cascades, the CaM dependent kinases activate phospholipases (e.g. phospholipase C) that in return activate other protein kinases such as protein kinase C.

CAMP

10

15

20

25

30

35

The cyclic AMP is produced by the enzyme adenylate cyclase in response to extracellular signals. Certain G-proteins stimulate the activity of adenylate cyclase which converts ATP to cAMP and PPi. Two molecules of cAMP bind to each of two regulatory subunits of cAMP dependent protein kinase which in turn dissociate from the two catalytic subunits of the heterotetramer R_2C_2 . Upon release of the C-subunits, they become active and phosphorylate substrate proteins at Ser and Thr residues. The process leading from binding of extracellular

molecules to their receptors, the transmission of the stimuli into the cell, the activation of adenylate cyclase and the subsequent activation of cAMP dependent protein kinase is one of two major signal transduction pathways in eukaryotic cells. Since the phosphorylation of proteins is a posttranslational modification of proteins, the kinases are described in the class "signal transduction."

SARA

10

15

20

25

30

35

Members of the transforming growth factor (TGFB) superfamily signal through a family of cell-surface transmembrane serine/threonine kinases, known as type I and type II receptors (Heldin et al., 1997; Attisano and Wrana, 1998; Kretzschmar and Massagué, 1998). Ligand induces formation of heteromeric complexes of these receptors, and signaling is initiated when receptor. I is phosphorylated and activated by the constitutively active kinase of receptor II (Wrana et al., 1994). The activated type I receptor kinase then propagates the signal to a family of intracellular signaling mediators known as Smads (contraction of the C-elegans Sma and Drosophila Mad genes which were the first identified members of this class of signaling effectors).

Three classes of Smads with distinct functions have been defined: the receptor-regulated Smads, which include Smadl, 2, 3, 5, and &; the common mediator Smad, Smad4; and the antagonistic Smads, which include Smadb and 7 (Heldin et al., 1997; Attisano and Wrana, 1998 ; Kretzschmar and Massagué, 1998). Receptorregulated Smads (R-Smads) act as direct substrates of specific type I receptors, and the proteins are phosphorylated on the last two serines at the carboxyl terminus within a highly conserved SSXS motif (Macías-Silva et al., 1996 ; Abdollah et al., 1997 ; Kretzschmar et al., 1997 ; Liu et al., 1997b ; Souchelnytskyi et al., 1997). Regulation of R-Smads by the receptor kinase provides an important level of specificity in this system. Thus, Smad2 and Smad3 are substrates of TGFR or activin receptors and mediate signaling by these ligands (Macías-Silva et al., 1996; Liu et al., 1997b ; Nakao et al., 1997), whereas Smadl, 5, and 8 are targets of BMP receptors and propagate BMP signals (Hoodless et al., 1996 ; Chen et al., 1997b ; Kretzschmar et al., 1997 ; Nishimura et al., 1998). Once phosphorylated, R-Smads associate with the common Smad Smad4 (Lagna et al., 1996 & Zhang et al.,

1997), and mediate nuclear translocation of the heteromeric complex. In the nucleus, Smad complexes then activate specific genes through cooperative interactions with DNA and other DNA-binding proteins such as FASTL, FASTZ, and Fos/Jun (Chen et al., 1996, Chen et al., 1997a; Liu et al., 1997a; Labbé et al., 1998; Zhang et al., 1998; Zhou et al., 1998). In contrast to R-Smads and Smad4, the antagonistic Smads, Smad6 and 7, appear to function by blocking ligand-dependent signaling (reviewed in Heldin et al., 1997).

10

15

20

25

30

35

Phosphorylation of R-Smads by the type I receptor essential for activating the TGFB signaling pathway (Heldin et al., 1997 ; Attisano and Wrana, 1998 ; Kretzschmar and Massaqué, 1998). However, little is known of how Smad interaction with receptors is controlled. A novel Smad2/Smad3 interacting protein has been described (Tsukazaki T. et al., 1998) that contains a double zinc finger, or FYVE domain, and which has been called SARA (Smad anchor for receptor activation). The SARA motif recruits Smad2 into distinct subcellular domains and co-localizes interacts with TGFR receptors. TGFR signaling induces dissociation of Smad2 from SARA with concomitant formation of Pbem2/Smad4 complexes and nuclear translocation. deletion of the FYVE domain in SARA causes mislocalization of Smad2 and inhibits TGFR-dependent transcriptional responses. Thus, SARA defines a component of TGFB signaling that functions to recruit Smad2 to the receptor by controlling the subcellular localization of Smad.

References: Abdollah et al. (1997) J. Biol. Chem. 272, 27678-27685; Attisano et al. (1998) Curr. Opin. Cell Biol. 10, 168-194; Chen et al. (1996) Nature 383, 691-696; Chen et al. (1997a) Nature 389, 85-89; Chen et al. (1997b) Proc. Natl. Acad. Sci. USA 94, 12938-12943; Heldin et al. (1997) Nature 390, 465-471; Hoodless et al. (1996) Cell 85, 489-500; Kretzschmar et al. (1998) Curr. Opin. Genet. Dev. 8, 103-111; Kretzschmar et al. (1997) Genes Dev. 11, 984-995; Labbé et al. (1998) Mol. Cell 2, 109-120; Lagna et al. (1996) Nature 383, 832-836; Liu et al. (1997a) Genes Dev. 11, 3157-3167; Liu et al. (1997b) Proc. Natl. Acad. Sci. USA 94, 10669-10764; Macías-Silva et al. (1996) Cell 87, 1215-1224; Nakao et al. (1997) EMBO J. 16, 5353-5362; Nishimura et al. (1998) J. Biol. Chem.

273, 1872-1879; Souchelnytskyi et al. (1997) J. Biol. Chem. 272, 28107-28115; Tsukazaki et al. (1998) Cell 95, 779-791; Wrana et al. (1994) Nature 370, 341-347; Zhang et al. (1997) Curr. Biol. 7, 270-276; Zhang et al. (1998) Nature 394, 909-913; Zhou et al. (1998) Mol. Cell 2, 121-127.

Calcium

5

10

The bivalent cation Ca²⁺ is, along with cAMP, one of the two major second messengers in eukaryotic cells. Its intracellular concentration is tightly regulated and usually kept very low compared to the cell's environment. Ca2+ binding proteins and transporters (Gap junction, Voltage-gated, second messengergated) help to sequester huge amounts of the ion in various organelles from where Ca2+ can be released upon extracellular 15 stimuli. E.g. the contraction of the muscle is dependent on the presence of Ca²⁺ ions which are readily transported back into the organelles in order for the muscle to relax. transduction, Ca²⁺ functions as a second messenger that activates Ca2+ dependent processes through the activation of Ca2+/calmodulin 20 dependent protein kinases (CaM kinases) which are the major. effector molecules of Ca2+. In the signaling cascades, the CaM dependent kinases activate phospholipases (e.g. phospholipase () that in return activate other protein kinases such as protein kinase C-

25 Rab proteins

30

35

In eukaryotic cells the compartmentalization of processes is tight regulation of processes a prerequisite for a activities. The cells contain a highly dynamic set of membrane compartments that are responsible for packaging, secreting, and recycling proteins and other Trafficking between organelles within the secretory pathway occurs as vesicles derived from a donor compartment fuse with specific acceptor membranes, resulting in the directional transfer of cargo molecules. This process is tightly controlled by the Rab/Ypt family of proteins (reviewed by Novick and Zerial, 1997), a branch of the superfamily of small GTPases. Rab proteins regulate a variety of functions, including vesicle translocation and docking at specific fusion sites. Rabs may also play critical roles in higher order processes such as modulating

the levels of neurotransmitter release in neurons, a likely mechanism in synaptic plasticity that underlies learning and memory (Geppert and SUdhof, 1998).

Small GTPases share a common three-dimensional fold that, in the GTP bound state, can bind a variety of downstream effector proteins. GTP hydrolysis leads to a conformational change in the "switch" regions that renders the GTPase unrecognizable to its effectors. In this way, by localizing and activating a select set of effectors, a common structural motif is used to control a wide array of distinct cellular processes.

5

15

20

25

30

35

The final steps in membrane fusion are likely to be driven by a set of proteins known as SNAREs. After a vesicle becomes the cytoplasmic domains of VAMP (also synaptobrevin) and syntaxin on opposing membranes, in combination with a SNAP-25 molecule, coalesce into an elongated -helical bundle (Poirier et al., 1998 ; Sutton et al., 1998), which may lead to fusion. Because numerous SNARE isoforms have been identified that localize to distinct membrane compartments, it was originally proposed that the specificity of interaction between the SNARE proteins accounted for the specificity in: membrane trafficking. Recent results, however, suggest that SNAREs are not specific in their ability to form complexes in suggesting that trafficking specificity additional factors (Yang et al., 1999). In this regard, Rab proteins are strong candidates for governing the specificity of vesicle trafficking. Like the SNAREs, many isoforms (40) of the Rab family have been identified that localize to specific membrane compartments (reviewed by Novick and Zerial, 1997).

Concomitant with the SNARE cycle, Rab proteins undergo a intricate cycle of membrane and protein interactions. Rabs are posttranslationally modified at C-terminal cysteines by the addition of two geranylgeranyl groups, which mediate membrane association when the Rab is in the GTP-bound state. After guanine nucleotide hydrolysis occurs, the Rab is extracted from the membrane upon forming a complex with a cytosolic GDP-dissociation inhibitor (GDI). This cytosolic intermediate is then recycled onto a newly forming vesicle, most likely through a secondary factor termed a GDI dissociation factor (GDF), which displaces GDI. After the Rab becomes membrane bound, a guanidine nucleotide

exchange factor (GEF) promotes release of GDP and the subsequent loading of GTP. In its GTP-bound conformation, the Rab is then free to associate with its specific set of effectors, which can in turn trigger events leading to the eventual fusion of the vesicle with a target membrane. To complete the cycle, perhaps after or concurrent with membrane fusion, a GTPase activating protein (GAP) accelerates nucleotide hydrolysis, switching off the GTPase. The remaining GDP-bound Rab can then participate in a new round of fusion

10

- 15

20

25

35

Rab interactions with effectors are likely to regulate vesicle targeting and membrane fusion in three ways. First, a Rab may specifically facilitate vectorial vesicle transport. Vesicles transported from their of site origin to. acceptor compartments likely through associations with cytoskeletal elements and transport motors. A protein has been identified with domain structure that suggests a connection between cytoskeleton and the Rabs. This protein, called Rabkinesin-b, contains a kinesin-like ATPase motor domain followed by a coiledcoil stalk region and a RBD that specifically binds Rabb (Echard et al., 1998). An additional link with the cytoskeleton is provided by the Rab effector, Rabphilin-3A. Rabphilin-3A has been shown in vitro to interact with -actinin, an actin-bundling protein, but only when not bound to Rab3A (Kato et al., 1996). These results raise the intriguing possibility that Rab proteins regulate vesicle interactions with the cytoskeleton and thereby play an active role in targeting vesicles to their appropriate destinations.

Second Rab proteins may regulate membrane trafficking at the vesicle docking step. A number of Rab effectors, including 30 Rabaptin-5, EEAl, Rabphilin-3A, and Rim, may serve as molecular tethers. Each effector protein contains a RBD, followed by a linker region (some having the potential to form elongated coiled-coil structures), and a domain capable of interacting with a second Rab or the target membrane. Rabaptin-5, for example, contains two RBDs, one near the N terminus that specifically recognizes Rab4 and a second near the C terminus that binds Rab5 (Vitale et al., 1998). Both Rim, which is localized to the target membrane, and Rabphilin-3A, which is localized to the vesicle contain N-terminal RBDs and C-terminal Ca2+-binding C2

domains, implicating these effectors in synaptic vesicle localization or docking in response to Ca2+ influx (Wang et al., 1997). Tethering effectors may also recognize protein complexes on the acceptor membrane. Sec4p, a yeast Rab3A homolog, interacts with the exocyst (Guo et al., 1999), a complex of seven or more subunits that is assembled at sites of vesicle fusion along the plasma membrane. The exocyst complex may therefore function as a landmark for Rab/effector-mediated vesicle docking.

Third, once a vesicle has become tethered to its fusion site, Rab proteins may selectively activate the SNARE fusion machinery. The mechanism of this activation is unknown but may involve direct interactions of Rabs or, more likely, their effectors with SNAREs. For example, Hrs-2 is a protein that binds to SNAP-25 and contains a Zn2+-finger motif characteristic of -Rab-binding proteins such as Rabphilin-3A, Rim, EEAL, and Noc2, --suggesting that Hrs-2 may form a physical link between Rabs and . SNAREs (Bean et al., 1997). In addition, certain mutations in the syntaxin-binding protein Slylp, the Seclp homolog utilized in ER to Golgi trafficking, eliminate the requirement for Yptlp, a Rab protein that functions at this trafficking step (Dascher et al., 1991). Rabs may therefore regulate SNARE associations through Secl family members. In support of this idea, a Rab effector was recently found to interact with a vacuole Rab, a Seclp homolog, and a SNARE protein (Peterson et al., 1999), which suggests that this effector serves to connect Rab and SNARE function. In this way, Rabs and their effectors may facilitate the correct pairing of SNAREs.

References: Dascher et al. (1991). Mol. Cell. Biol. 11, 872-885; Echard et al. (1998). Science. 279, 580-585; Geppert et al. (1998). Annu. Rev. Neurosci. 21, 75-95; Guoet al. (1999). EMBO J. 18, 1071-1080; Kato et al. (1996). J. Biol. Chem. 271, 31775-31778; Novick et al. (1997). Curr. Opin. Cell Biol. 9, 496-504; Peterson et al. (1999). Curr. Biol. 9, 159-162; Poirier et al. (1998). Nat. Struct. Biol. 5, 765-769; Vitale et al. (1998). EMBO J. 17, 1941-1951; Wang et al. (1997). Nature. 388, 593-598; Yang et al. (1999). J. Biol. Chem. 274, 5649-5653.

<u>Kinases</u>

10

15

20

25

30

35

Reversible posttranslational modifications of proteins are major means of regulating cellular activities. Among the various modifications that are carried out by the cells, the addition of phosphoryl groups to Ser/Thr or Tyr residues is the most important and widely used. The phosphorylation of proteins is accomplished by protein kinases, while the reverse reaction, the removal of phosphoryl groups, is carried out by phosphatases. Kinases / Phosphatases regulate key positions e.g. in the processes of cell proliferation, differentiation and communication/signaling. These processes must be tightly 10 regulated in order to maintain a steady state level of cellular fate. Mis-regulation of kinase activities (or that of phosphatases) is made responsible for a multitude of disease processes such as oncogenesis, inflammatory processes, arteriosclerosis, and psoriasis. 15

Protein kinases constitute the largest protein family that is currently known. Several hundred kinases have been identified already. Classically, kinases are subdivided into two classes based on the amino acid residues in their substrates that are phosphorylated by the particular enzymes. The kinases specifically add phosphoryl groups from adenosine triphosphate (ATP) or, less frequently, guanosine triphosphate (GTP), either to serine and/or threonine or to tyrosine residues of substrate proteins. An estimated 1,000 to 10,000 proteins present in a typical mammalian cell are believed to be regulated also by the action of protein kinases.

- 20

25

35

Protein kinases are frequently integral parts of signaling cascades that transmit extracellular stimuli (e.g. hormones, neurotransmitters, growth— or differentiation factors) into the cell and result in various responses by the cells. The kinases play key roles in these cascades as they constitute a sort of 'molecular switches' turning on or off the activities of other enzymes and proteins, e.g. metabolic, regulatory, channels and pumps, receptors, cytoskeletal, transcription factors.

The regulation of kinase activities is accomplished by various means:

The best characterized example for the regulation via regulatory subunits is the cAMP-dependent protein kinase (PKA) which is also a prototype for second messenger activated protein

kinases. This enzyme consists of a heterotetramer of two catalytic (C) and two regulatory (R) subunits. Upon binding of two molecules of second messenger (cAMP) in each R subunit, the catalytic subunits are released and active. Both of the catalytic and the regulatory subunits several isoforms exist. The combination of catalytic and regulatory subunits determines the localization of the holoenzyme and also the substrate spectrum that is available for phosphorylation. The consensus pattern necessary to be present in the substrate for PKA action is RRXS/T where X can be any amino acid.

The casein kinase II comprises another examples for holoenzymes that consist of catalytic and regulatory subunits. Other kinases that are activated by second messengers are cGMP-dependent protein kinase and Protein kinase C (PKC) which is activated by diacylglycerol, which in turn is produced by phospholipases by cleavage of phosphatidylcholine.

10

15

20

25

30

35

Receptor kinases usually consists of an extracellular domain which can bind effector molecules (e.g. growth factors and hormones) and transfer the stimulus to the intracellular domain of these proteins which usually is a protein tyrosine kinase. Other tyrosine kinases lack an extracellular domain but are associated with receptors which transfer the signal after effector binding by activating the associated protein kinase enzyme (e.g. Src kinase family; Src, Blk, Fgr, Fyn, Lck Lyn, Yes and Janus kinase family; Jakl-3, Tyk2).

Dysfunction of kinases, e.g. caused by non-functioning regulation, can be the cause of inflammatory diseases and uncontrolled proliferation. v-Src which is a truncated version of the C-Src protooncogene tyrosine kinase is a classical example for this process as v-Src does not contain the regulatory domain of the cellular gene and is thus constitutively active.

Several categories of proteins are coded for by clones of the invention within the overall group of "Signal transduction"and include, among others, the following:

<u>Discs-large family:</u> In Drosophila more than 50 genes are discribed in which mutation leads to loss of cell proliferation control indicating that they are tumor suppressor genes. Most of

these genes have mammalian homologs. The Drosophila 'discs large' tumor suppressor protein, Dlg, is the prototype of a family of proteins termed MAGUKs (membrane-associated guanylate kinase homologs). MAGUKs are localized at the membrane-cytoskeleton interface, usually at cell-cell junction, where they appear to have both structural and signaling roles. They contain several distinct domains, including a modified guanylate kinase domain, an SH3 motif, and 1 or 3 copies of the DHR (GLGF/PDZ) domain. Recessive lethal mutations in the 'discs large' tumor suppressor gene interfere with the formation of septate junctions (thought to be the arthropod equivalent of tight junctions) between epithelial cells, and they also cause neoplastic overgrowth of imaginal discs, suggesting a role for cell junctions in proliferation control. These proteins can find application in modulating/blocking the guanylate cyclase-pathway. Clones in this category include: amy2_12d7.

10

15

35

· Proteins with a WW Domain: Proteins that contain a WW domain which has been originally described as a short conserved 20 region in a number of unrelated proteins, among them dystrophin, the gene responsible for Duchenne muscular dystrophy. The domain, which spans about 35 residues, is repeated up to 4 times in some proteins. It has been shown to bind proteins with particular proline-motifs; TAPI-P-P-EAPI-Y, and thus resembles somewhat SH3 25 domains. This domain is frequently associated with other domains typical for proteins in signal transduction processes. Examples of proteins containing the WW domain are Dystrophin, Utrophin, vertebrate YAP protein (binds the SH3 domain of the Yes oncoprotein), murine NEDD-4 (embryonic development and 30 differentiation of the central nervous system), IQGAP (human GTPase activating protein acting on ras). Therefore these proteins should be involved in intracellular signal transduction. Diseases associated (as potentially diagnostic, therapeutic, causative, and/or related, etc...) with these proteins include as reported by OMIN 1) Muscular Dystrophy, Pseudohypertrophic Progressive Duchenne and Becker Types (OMIN *310200). Clones in this category include: tes3_lld2l.

<u>Ion-Transporters:</u> For signalling stringent control od ion fluxes over biological membranes is of the essence. Several trans-membrane ion-chennel-proteins key elements of signal transduction pathways. Clones in this category include: amy2_10p7 and amy2_2f18.

RING-finger proteins: A Zinc finger motif of the C3HC4 type (the so-called RING finger domain) is involved in mediating protein-protein interactions. Proteins containing a RING-finger are: mammalian V(D)J recombination activating protein (RAGL), mouse rpt-l, human rfp, human 52 Kd Ro/SS-A protein and others. The family of RING finger proteins contains a number of oncogenes. For example PML, a probable transcription factor, BRCAL, the mammalian cbl- and bmi-l proto-oncogenes. Clones in this category include: amy2_10h17.

Phosphatases: Proper targeting of PTPs is essential for many cellular signalling events including antigen induced proliferative responses of B and T cells. The physiological significance of PTPs is further unveiled through mice gene knockout studies and human genome sequencing and mapping projects. Several PTPs are shown to be critical in the pathogenesis of human diseases, as shown by over 290 entries in OMIN. Clones in this category include: tes3_3lj20.

Phosphoproteins: Some paraneoplastic syndromes affecting the nervous system are associated with antibodies that react with neuronal proteins and the causal tumor (onconeuronal antigens). Several of these antibodies are markers of specific neurologic syndromes associated with distinct types of cancer. One of the antigenes recognised by such antibodies is Ma-l, the neuron- and testis-specific protein l. The expression of Mal mRNA is highly restricted to the brain and testis. Subsequent analysis suggested that Mal is likely to be a phosphoprotein (see OMIN *604010). Clones in this category include: tes3_5k22.

Transmembrane proteins

10

25

30

35

Membrane region prediction was effected using the ALOM2 software (Klein et al., 1985; version 2 by K. Nakai). Similar to

many other methods, the Kyte & Doolitle (1982) amino acid hydrophobicity scale is used in ALOM2 as the primary variable for classifying sequences in terms of their localization. High prediction accuracy is achieved through the system of intelligent decision rules and the utilization of a carefully selected training data set. The method also generates reliability estimates which makes it possible to distinguish between membrane-spanning proteins (I, intrinsic) and globular proteins with regions of high hydrophobicity buried in the core.

10 For a protein of length L_1 the block of length 1 with maximum hydrophobicity is found:

$$\max H = \max(1/l) \sum_{\substack{i=k\\k=1,\dots,L-l+1}}^{k+l-1} H_i$$

where \mathcal{H}_i represents the hydrophobicity of an individual residue.

Let P(I/maxH) and P(E/maxH) be the conditional probabilities that a protein is integral or peripheral, respectively, given its value of maximal hydrophobicity maxH, and let P(I) and P(E) be the prior probabilities of intrinsic and extrinsic membrane proteins estimated from the training set. Then a sequence is assigned to E if

P(E/maxH) > P(I/maxH)

30

or, after applying the Bayes rule,

P(E)P(maxH/E) > P(I)P(maxH/I)

where the conditional probabilities P(maxH/E) and P(maxH/I)

25 can be determined based on the estimates of probability
distributions of maxH in both groups.

Discriminant analysis allows to simplify this task by calculating the odds P(E/MaxH):P(I/maxH) as e^b , where b is the left-hand side of a linear or quadratic inequality. For example, for the window of length 17, the protein is allocated to the

peripheral category E based on the empirically derived quadratic inequality:

1.05(maxH)²+12.30maxH+17.49 >0,

whereas the optimal inequality for assigning membrane 5 proteins (category I) is linear:

-9.02maxH + 14.27 > 0

The odds parameter can be made more or less stringent. For example, one can require odds at least 1:10 for a protein to be classified as integral. This leads to higher selectivity but less sensitivity.

The boundaries of membrane-spanning regions in putative membrane proteins are detected by means of an iterative procedure whereby the most hydrophobic region corresponding to the value maxH is considered to be membrane and removed from the sequence. The classification procedure is then repeated again for the remaining sequence, and, if such a protein is again classified as integral, the next most hydrophobic region is considered.

Reference: Klein, P., Kanehisa, M., DeLisi, C. (1985) The detection and classification of membrane-spanning proteins. Biochem Biophys Acta 815: 468-476

Transcription factors

10

15

20

25

30

Purified eukaryotic RNA polymerase II is unable to initiate promoter-specific transcription. A family of factors that collectively confer RNAPII promoter specificity is known as the general transcription factors (GTFs). They include the TATA-binding Protein (TBP) TFIIB, TFIIE, TFIIF and TFI IH. These factors are conserved among all eukaryotes.

RNAPII complexes containing the entire set of GTFs or a subset of GTFs together with other proteins have been isolated from mammalian and yeast cells. Although purified RNAPII and GTFs are sufficient for promoter-specific initiation, this system fails to respond to activators. This is mediated by a further complex termed mediator complex which associates with the

WO 01/98454 PCT/IB01/02050 carboxy-terminal heptapeptide domain (CTD) of the largest subunit of RNAPII.

Purification of human RNAPII complexes resulted in two distinct forms of human RNAPII after analysis of functional properties. One complex contained chromatin remodeling activities but was devoid of GTFs. The other complex did not contain factors that modify chromatin but contained a subset of SRB/mediator subunits and GTFs and other polypeptides that mediate transcriptional activation, a scenario similar to that reported for yeast.

A complex designated NAT (~20 SU) for negative regulator of transcription contains RNAPII, Cdk8, homologs of the yeast mediator complex as well as Rgrl and Srbl0/ll known as negative regulators of transcription.

A complex with striking similar structural and functional properties to NAT has been identified designated SMCC (~15 SU) (SRB/mediator coactivator complex), that can also mediate transcriptional activation.

The SMCC complex includes all reported NAT subunits

20 including subunits of the TRAP complex. TRAP is a coactivator complex isolated on the basis of its interaction with the thyroid hormone receptor. Another coactivator complex DRIP, isolated on the basis of its ability to interact with the vitamin D3 receptor, contains novel subunits as well as subunits of NAT/SMCC and TRAP complexes.

The effects of each of these coactivator complexes is dependent on the TFIID complex. It is not known if the T AF subunits of TFIID are required. It is likely that new coactivator complexes will be uncovered containing both novel and previously defined components.

30

35

Beside the huge amount of transcription factors which can be part of the RNAIIP holoenzyme or the coactivator complexes there is an even larger quantity of specific transcription factors binding to promoter elements within the DNA sequences of a given gene leading to activation or repression of transcription. A

broad range of cellular responses like differentiation, proliferation, cell death and others are elicited through activating or repressing the transcription of target genes.

There are at least five superclasses of transcription 5 factors:

1. Superclass contains members with characteristic basic domains:

Members are:

Leucine zipper factors, where the basic domain is followed 10 by a leucine zipper of repeated leucine residues at every seventh position. The zipper mediates protein dimerization as a prerequisite for DNA-binding.

Helix-loop-helix factors (bHLH) contain a DNA-binding basic region followed by a motif of two potential amphipathic alphahelices connected by a loop of variable length also mediating dimerization.

Factors with a combination of Helix-loop-helix and leucine zipper-

Further members of this superclass are NF-l, RF-X, and bHSH 20 like proteins.

2. Superclass comprises factors containing zinc-coordinating DNA-binding domains.

Members are:

Proteins with Cys4 zinc finger of nuclear receptor type.

25 where two such motifs differing in size, composition and function are present in each receptor molecule. Each finger comprises 4 cysteine residues coordinating one zinc ion. The second half including the second cysteine pair has alpha-helix conformation and the helix of the first finger binds to the DNA through the major groove. The sequence between the first two cysteines of the second finger mediates dimerization upon DNA-binding. This class includes the steroid hormone receptors and the thyroid hormone

receptor-like factors. Other diverse cys4 zinc fingers have a motif of GATA-type.

Proteins with Cys2His2 zinc finger domain(s). Each finger comprises 2 cysteine and 2 histidine residues coordinating one zinc ion, and in some cases one histidine is replaced by another cysteine. The zinc ion is essential for DNA-binding.

Proteins with Cysb cysteine-zinc cluster(s). Six cysteine residues coordinate two zinc ions, i. e. two of the thiol groups are coordinating two zinc ions each. Present in many fungal regulators.

Zinc fingers of alternating composition.

3- Superclass contains factors of helix-turn-helix type.

Members are:

10

20

Proteins with homeo domains. Homeo domains are three

15 consecutive alpha-helix structures. Helix 3 contacts mainly the major groove of the DNA, some contacts at the minor groove are observed as well. Helix 2 and 3 resemble the helix-turn-helix structure of prokaryotic regulators.

Proteins with Paired box domain(s). This is a DNA-binding domain of approximately 130 amino acid residues. Its N-terminal half is basic, its C-terminal half is highly charged in general. It probably comprises 3 alpha-helices.

Proteins with Fork head / winged helix domain(s). This domain was identified by homology between HNF-3A and fkh. The domain comprises approx. 110 AA. Analysis of the crystal structure has revealed a compact structure of three alphahelices, the third alphahelix being exposed towards the major groove of the DNA. The domain also exerts minor groove contacts. Upon binding to DNA, it induces a bend of 13 degree.

30 Heat shock factors

Proteins with Tryptophan clusters. The tryptophan clusters comprise several tryptophan residues with a spacing of 12-21

amino acid residues; the subclass of myb-type DNA-binding domains typically exhibit a spacing of 19-21 amino acid residues.

Proteins with TEA domain(s). The TEA domain has been identified as a region which is conserved among the transcription factors TEF-1. TECl and abaA. This domain in TEF-1 has been shown to interact with DNA, although two additional regions may also contribute to DNA-binding. It is predicted to fold into three alpha-helices, with a randomly coiled region of 16-18 amino acid residues between helices 1 and 2, and a short stretch between helices 2 and 3 of 3-8 residues.

4- Superclass contains beta-Scaffold Factors with Minor Groove Contacts

Members are:

10

Proteins with RHR (Rel homology) region.

15 The structure of the Rel-type DBD exhibits a bipartite subdomain structure, each subdomain comprising a beta-barrel with five loops that form an extensive contact surface to the major groove of the DNA. Particularly, the first loop of the N-terminal subdomain (the highly conserved recognition loop) performs 20 contacts with the recognition element on the DNA, but other loops are involved. The fact that the main DNA-contacts are made through loops has been suggested to provide a high degree of flexibility in binding to a range of different target sequences. Augmenting interactions are achieved by two alpha-helices within the N-terminal Part that form strong minor groove contacts to the 25 A/T-rich center of the B-element. In pb5, the sequence between both alpha-helices is much shorter and even helix 2 is truncated. The second, C-terminal domain is necessary mainly for protein dimerization.

30 p53 proteins

MADS (MCM1-agamous-deficiens-SRF) box proteins. Proteins of this class comprise a region of homology. The DNA-binding domain also comprises the dimerization capability. In the DNA-bound dimer (shown for SRF), two antiparallel amphipathic alpha-helices

(alpha-I), form a coiled coil and are oriented approximately parallel on the minor groove. These helices make minor and major groove contacts, the N-terminal extensions form minor groove contacts. The bound DNA is bent and wrapped around the proteinly texhibits a compressed minor groove in the center and widened minor groove in the flanks.

Beta-Barrel alpha-helix transcription factors.

TATA-binding proteins

HMG proteins

10 Proteins of this class comprise a region of homology with the chromosomal non-histone HMG proteins such as HMGL. This region comprises the DNA-binding domain which in some instances such as HMGL mediates sequence-unspecific in other cases such LEF-L sequence-specific binding to DNA. This domain exhibits a typical L-shaped conformation made up of 3 alpha-helices and an extended N-terminal extension of the first helix. The latter together with helix L which contains a kink form the long arm of the L whereas helices L and 2 form the short arm. Binding to the minor groove induces a sharp bending of the DNA by more than 90 degree away from the bound protein. The overall topology of the DNA-protein complexes resembles somewhat that of the TBP-TATA box complex.

Heteromeric CCAAT factors

Proteins with Grainyhead domain(s)

Cold-shock domain factors. Cold-shock domain proteins are characterized by a highly conserved region first found in prokaryotic cold-shock proteins. This domain is a single-stranded nucleic acid-binding structure interacting with DNA or RNA. It consists of an antiparallel five-stranded beta-barrel, the strands of which are connected by turns and loops. Within this structure, a three-stranded beta-strand contains a conserved RNA-binding motif, RNP1. Not all CSD proteins are transcription factors. Those which specifically bind to a certain sequence are termed Y-box proteins. Proteins of this class were previously

called protamine-like domain proteins because of having a highly positively charged domain with interspersed proline residues.

Proteins with Runt homology domain

The members of this transcription factor class have been

5 identified on the basis of their homology to a defined region
within the Drosophilia protein Runt. The runt domain is part of
the DNA-binding domain of these factors. It consists mainly of
beta-strands, does not contain alpha-helical regions and seems to
be most similar to the palm domain found in DNA polymerase beta

10 (rat).

5. Superclass contains other transcription factors like Copper fist proteins, HMGI(Y), STAT, Pocket domain proteins and Ap2/EREBP-related factors.

The classification of transcription factors originates from 15 TRANSFAC database:

http://transfac.gbf.de/TRANSFAC/

Reference: Heinemeyer

20

25

30

Several categories of proteins are coded for by clones of the invention within the overall group of "Transcription Factors" and include, among others, the following:

Homeobox-proteins: Homeodomain-containing transcription factors are essential for a variety of processes in vertebrate development, including organogenesis. They have been shown to regulate cell proliferation, pattern segmental identity anddetermine cell fate decisions during embryogenesis. For example, In zebrafish emx2 mRNAs are found in the dorsal telencephalon, parts of the diencephalon and the otocyst. The human homologue Emx2 appears to be already expressed in 8.5 day embryos. It is also expressed in the presumptive cerebral cortex, olfactory bulbs, in some neuroectodermal areas in embryonic head including olfactory placodes in earlier stages and olfactory epithelia later in development. Mutants of the D. melanogaster gene "mempty spiracles" display spiracles devoid of filzkorper,

no antenna and an open head. Clones in this category include: amy2_14m16.

Proteins with myc-type, helix-loop-helix dimerization domain signature(s). This helix-loop-helix domain mediates protein dimerization has been found in various multimeric transcrpition factors. Clones in this category include: tes3_l8nl4.

<u>Transcriptional silencers:</u> In addition to transcription factors, other proteins, such as YDL153c of Saccharomyces cerevisia are responsible for silencing of genes. Clones in this category include: amy2_2f22.

10

15

20

25

30

35

Proteins regulating transcription factors: The activity of several transcription factor is regulated by the binding or dissociation of other proteins or by phosphorylation or dephosphorylation of the transcription factor. For example, I-kappa-B-related protein interacts with the transcription factor NF-kB. I-kappa-B-alpha mutations contribute to constitutive NF-kappaB activity in cultured and primary HRS (Hodgkin/Reed-Sternberg) cells and are therefore involved in the pathogenesis of Hodgkin's disease (HD) patients. Clones in this category include: amy2_lcl2.

Signal transducing proteins: Beta-transducin subunits of G-proteins contain WD-4D repeats. The beta subunits seem to be required for the replacement of GDP by GTP as well as for membrane anchoring and receptor recognition. Due to the zinc finger the novel protein seems to be a new molecule involved in signal transduction and transcription. These proteins have been reported by OMIN to be associated (as potentially diagnosticated by OMIN to be associated (as potentially diagnosticated diseases: 1) essential hypertension (OMIN *139130). Clones in this category include: tes3_llc22.

* * *

The invention, therefore, specifically contemplates the following assemblages of materials, which track the above-identified fourteen functional groupings, that are useful in practicing the profiling aspects of the invention. One type of assemblage is nucleic acid-based and can include the following groupings of sequences and their derivatives: all sequences; human fetal

kidney library sequences; kidney derived sequences; human mammary carcinoma library sequences; mammary carcinoma derived sequences; human testis library sequences; testes derived sequences; cell cycle genes; cell structure and motility genes; differentiation and development genes; intracellular transport and trafficking genes; metabolism genes; nucleic acid management genes; signal transduction genes; transmembrane protein genes; and transcription factor genes. Other assemblages contain proteins or their corresponding antibodies or antibody fragments, divided along the same groupings.

Database Applications

10

15

20

25

35

Because they are human genes and gene products: the inventive molecules are useful as members of a database. Such a database may be used: for example: in drug discovery and rationale drug design or in testing the novelty and non-obviousness of newly sequenced materials. In addition: they are particularly suited in designing variants for the profiling (and other) applications described herein. Hence: the following discussion of electronic embodiments applies equally to such variants: which: naturally: will be generated and stored using a computer using known methodologies.

Accordingly, one aspect of the invention contemplates a database of at least one of the inventive sequences stored on computer readable media. Again, the individual sequences may be grouped with regard to the individual functional and structural groups mentioned above. While the individual sequences of a database may exist in printed form, they are preferably in electronic form, as in an ascii or a text file. They may also exist as word processing files or they may be stored in database applications like DB2. Sybase, Oracle, GCG and GenBank. One skilled in the art will understand the range of applications suitable for using and storing the electronic embodiments of the invention.

"Computer readable media" refers to any medium which can be read and accessed by a computer. These include: magnetic storage media, like floppy discs, hard drives and magnetic tape; optical storage media, like CD-ROM; electrical storage media, like RAM and ROM; and hybrids of these categories, like magnetic/optical

storage media. One skilled in the art will readily understand the scope of computer readable media and how to implement them.

Biological Activities and Assays for Implementing Therapeutic and Diagnostic Applications

This section provides assays for biological activity that are useful in characterizing and quantifying the biological activity of the inventive molecules and their derivatives, which is relevant to the pharmacological effects of the inventive molecules. As used in this section, it will be understood that "protein" may also refer to the inventive antibodies (including fragments).

Cytokine and Cell Proliferation/Differentiation Activity

10

15

20

25

A protein of the present invention may exhibit cytokine; cell proliferation (either inducing or inhibiting) or cell differentiation (either inducing or inhibiting) activity or may induce production of other cytokines in certain cell populations. Many protein factors discovered to date; including all known cytokines; have exhibited activity in one or more factor dependent cell proliferation assays; and hence the assays serve as a convenient confirmation of cytokine activity. The activity of a protein of the present invention is evidenced by any one of a number of routine factor dependent cell proliferation assays for cell lines including; without limitation; 32D; DA2; DA1; 123; T10; B9; B9/11; BaF3; MC9/6; M + (preB M +); 2E8; RB5; DA1; 123; T1165; HT2; CTLL2; TF-1; Mo7e and CMK.

The activity of a protein of the invention may, among other means, be measured by the following methods:

Assays for T-cell or thymocyte proliferation include without limitation those described in: Current Protocols in Immunology.

30 Ed by J. E. Coligan, A. M. Kruisbeek, D. H. Margulies, E. M. Shevach, W. Strober, Pub. Greene Publishing Associates and Wiley-Interscience (Chapter 3, In Vitro assays for Mouse Lymphocyte Function 3.1-3.19; Chapter 7, Immunologic studies in Humans); Takai et al., J. Immunol. 137:3494-3500, 1986; Bertagnolli et al., Cellular Immunology 133:327-341, 1990; Bertagnolli et al., Cellular Immunology 133:327-341, 1991; Bertagnolli, et al., I. Immunol. 149:3778-3783, 1992; Bowman et al., I. Immunol. 152:1756-1761, 1994.

Assays for cytokine production and/or proliferation of spleen cells, lymph node cells or thymocytes include, without limitation, those described in: Polyclonal T cell stimulation, Kruisbeek, A. M. and Shevach, E. M. In Current Protocols in Immunology. J. E. e.a. Coligan eds. Vol 1 pp. 3.12.1-3.12.14, John Wiley and Sons, Toronto. 1994; and Measurement of mouse and human interleukin gamma, Schreiber, R. D. In Current Protocols in Immunology. J. E. e.a. Coligan eds. Vol 1 pp. b.8.1-b.8.8, John Wiley and Sons, Toronto. 1994.

Assays for proliferation and differentiation of 10 hematopoietic and lymphopoietic cells include, without limitation, those described in: Measurement of Human and Murine Interleukin 2 and Interleukin 4, Bottomly, K., Davis, L. S. and Lipsky, P. E. In Current Protocols in Immunology. J. E. e.a. 15 Coligan eds. Vol 1 pp. b.3.1-b.3.12. John Wiley and Sons. Toronto. 1991; devries et al., J. Exp. Med. 173:1205-1211, 1991; Moreau et al., Nature 336:690-692, 1988; Greenberger et al., Proc. Natl. Acad. Sci. U.S.A. 80:2931-2938, 1983; Measurement of mouse and human interleukin b-Nordan, R. In Current Protocols in Immunology. J. E. e.a. Coligan eds. Vol 1 pp. 6-6-1-6-5. John 20 Wiley and Sons, Toronto. 1991; Smith et al., Proc. Natl. Aced. Sci. U.S.A. 83:1857-1861, 1986; Measurement of human Interleukin 11-Bennett, F., Giannotti, J., Clark, S. C. and Turner, K. J. In Current Protocols in Immunology. J. E. e.a. Coligan eds. Vol 1 25 pp. 6.15.1 John Wiley and Sons, Toronto. 1991; Measurement of mouse and human Interleukin 9-Ciarletta, A., Giannotti, J., Clark, S. C. and Turner, K. J. In Current Protocols in. Immunology. J. E. e.a. Coligan eds. Vol 1 pp. 6.13.1, John Wiley and Sons, Toronto. 1991.

Assays for T-cell clone responses to antigens (which will identify, among others, proteins that affect APC-T cell interactions as well as direct T-cell effects by measuring proliferation and cytokine production) include, without limitation, those described in: Current Protocols in Immunology, Ed by J. E. Coligan, A. M. Kruisbeek, D. H. Margulies, E. M. Shevach, W Strober, Pub. Greene Publishing Associates and Wiley-Interscience (Chapter 3, In Vitro assays for Mouse Lymphocyte Function; Chapter b, Cytokines and their cellular receptors; Chapter 7, Immunologic studies in Humans); Weinberger et al.,

Proc. Natl. Acad. Sci. USA 77:6091-6095, 1980; Weinberger et al., Eur. J. Immun. 11:405-411, 1981; Takai et al., J. Immunol. 137:3494-3500, 1986; Takai et al., J. Immunol. 140:508-512, 1988.

Immune Stimulating or Suppressing Activity

5

10

15

20

25

30

35

A protein of the present invention may also exhibit immune stimulating or immune suppressing activity, including without limitation the activities for which assays are described herein. A protein may be useful in the treatment of various immune deficiencies and disorders (including severe combined immunodeficiency (SCID)), e-g., in regulating (up or down) growth and proliferation of T and/or B lymphocytes, as well as effecting the cytolytic activity of NK cells and other cell populations. These immune deficiencies may be genetic or be caused by vital (e-g-, HIV) as well as bacterial or fungal infections, or may result from autoimmune disorders. More specifically, infectious diseases causes by viral, bacterial, fungal or other infection may be treatable using a protein of the present invention, including infections by HIV, hepatitis viruses, herpesviruses, mycobacteria, Leishmania spp., malaria spp. and various fungal infections such as candidiasis. Of course, in this regard, a protein of the present invention may also be useful where a boost to the immune system generally may be desirable, i.e., in the treatment of cancer-

Autoimmune disorders which may be treated using a protein of the present invention include, for example, connective tissue disease, multiple sclerosis, systemic lupus erythematosus, rheumatoid arthritis, autoimmune pulmonary inflammation, Guillain-Barre syndrome, autoimmune thyroiditis, insulin dependent diabetes mellitis, myasthenia gravis, graft-versus-host disease and autoimmune inflammatory eye disease. Such a protein of the present invention may also to be useful in the treatment of allergic reactions and conditions, such as asthma (particularly allergic asthma) or other respiratory problems. Other conditions, in which immune suppression is desired (including, for example, organ transplantation), may also be treatable using a protein of the present invention.

Using the proteins of the invention it may also be possible to modify immune responses, in a number of ways. Down regulation

may be in the form of inhibiting or blocking an immune response already in progress or may involve preventing the induction of an immune response. The functions of activated T cells may be inhibited by suppressing T cell responses or by inducing specific tolerance in T cells, or both. Immunosuppression of T cell responses is generally an active, non-antigen-specific, process which requires continuous exposure of the T cells to the suppressive agent. Tolerance, which involves inducing non-responsiveness or anergy in T cells, is distinguishable from immunosuppression in that it is generally antigen-specific and persists after exposure to the tolerizing agent has ceased. Operationally, tolerance can be demonstrated by the lack of a T cell response upon reexposure to specific antigen in the absence of the tolerizing agent.

10

15

20

25

30

35

Down regulating or preventing one or more antigen functions (including without limitation B lymphocyte antigen functions (such as, for example, B7)), e.g., preventing high level lymphokine synthesis by activated T cells, will be useful in situations of tissue, skin and organ transplantation and in graft-versus-host disease (GVHD). For example, blockage of T cell function should result in reduced tissue destruction in tissue transplantation. Typically, in tissue transplants, rejection of the transplant is initiated through its recognition as foreign by T cells, followed by an immune reaction that destroys the transplant. The administration of a molecule which inhibits or blocks interaction of a B7 lymphocyte antigen with its natural. ligand(s) on immune cells (such as a soluble, monomeric form of a peptide having B7-2 activity alone or in conjunction with a monomeric form of a peptide having an activity of another B lymphocyte antigen (e.g., B7-1, B7-3) or blocking antibody), prior to transplantation can lead to the binding of the molecule to the natural ligand(s) on the immune cells without transmitting the corresponding costimulatory signal. Blocking B lymphocyte antigen function in this matter prevents cytokine synthesis by immune cells, such as T cells, and thus acts as an immunosuppressant. Moreover, the lack of costimulation may also be sufficient to anergize the T cells, thereby inducing tolerance in a subject. Induction of long-term tolerance by B lymphocyte antigen-blocking reagents may avoid the necessity of repeated

PCT/IB01/02050 WO 01/98454

administration of these blocking reagents. To achieve sufficient immunosuppression or tolerance in a subject, it may also be necessary to block the function of a combination of B lymphocyte antigens.

The efficacy of particular blocking reagents in preventing organ transplant rejection or GVHD can be assessed using animal models that are predictive of efficacy in humans. Examples of appropriate systems which can be used include allogeneic cardiac grafts in rats and xenogeneic pancreatic islet cell grafts in mice, both of which have been used to examine the immunosuppressive effects of CTLA4Ig fusion proteins in vivo as described in Lenschow et al., Science 257:789-792 (1992) and Turka et al., Proc. Natl. Acad. Sci USA, 89:11102-11105 (1992). In addition, murine models of GVHD (see Paul ed., Fundamental -15 Immunology: Raven Press: New York: 1989: pp. 846-847) can be used to determine the effect of blocking B lymphocyte antigen function in vivo on the development of that disease.

10

20

25

35

Blocking antigen function may also be therapeutically useful for treating autoimmune diseases. Many autoimmune disorders are the result of inappropriate activation of T cells that are reactive against self tissue and which promote the production of cytokines and autoantibodies involved in the pathology of the diseases. Preventing the activation of autoreactive T cells may reduce or eliminate disease symptoms. Administration of reagents which block costimulation of T cells by disrupting receptor:ligand interactions of B lymphocyte antigens can be used to inhibit T cell activation and prevent production of autoantibodies or T cell-derived cytokines which may be involved in the disease process. Additionally, blocking reagents may 30 induce antigen-specific tolerance of autoreactive T cells which could lead to long-term relief from the disease. The efficacy of blocking reagents in preventing or alleviating autoimmune disorders can be determined using a number of well-characterized animal models of human autoimmune diseases. Examples include murine experimental autoimmune encephalitis, systemic lupus erythmatosis in MRL/lpr/lpr mice or NZB hybrid mice, murine autoimmune collagen arthritis, diabetes mellitus in NOD mice and BB rats, and murine experimental myasthenia gravis (see Paul ed.,

Fundamental Immunology, Raven Press, New York, 1989, pp. 840-856).

Upregulation of an antigen function (preferably a B lymphocyte antigen function), as a means of up regulating immune responses, may also be useful in therapy. Upregulation of immune responses may be in the form of enhancing an existing immune response or eliciting an initial immune response. For example, enhancing an immune response through stimulating B lymphocyte antigen function may be useful in cases of viral infection. In addition, systemic viral diseases such as influenza, the common cold, and encephalitis might be alleviated by the administration of stimulatory forms of B lymphocyte antigens systemically.

10

15

20

25

30

35

Alternatively, anti-vital immune responses may be enhanced in an infected patient by removing T cells from the patient, costimulating the T cells in vitro with viral antigen-pulsed APCs either expressing a peptide of the present invention or together with a stimulatory form of a soluble peptide of the present invention and reintroducing the in vitro activated T cells into the patient. Another method of enhancing anti-viral immune responses would be to isolate infected cells from a patient, transfect them with a nucleic acid encoding a protein of the present invention as described herein such that the cells express all or a portion of the protein on their surface, and reintroduce the transfected cells into the patient. The infected cells would now be capable of delivering a costimulatory signal to, and thereby activate. T cells in vivo.

In another application, up regulation or enhancement of antigen function (preferably B lymphocyte antigen function) may be useful in the induction of tumor immunity. Tumor cells (e.g., sarcoma, melanoma, lymphoma, leukemia, neuroblastoma, carcinoma) transfected with a nucleic acid encoding at least one peptide of the present invention can be administered to a subject to overcome tumor-specific tolerance in the subject. If desired, the tumor cell can be transfected to express a combination of peptides. For example, tumor cells obtained from a patient can be transfected ex vivo with an expression vector directing the expression of a peptide having B7-2-like activity alone, or in conjunction with a peptide having B7-1-like activity and/or B7-3-like activity. The transfected tumor cells are returned to the

PCT/IB01/02050 WO 01/98454

patient to result in expression of the peptides on the surface of the transfected cell. Alternatively, gene therapy techniques can be used to target a tumor cell for transfection in vivo-

The presence of the peptide of the present invention having the activity of a B lymphocyte antigen(s) on the surface of the tumor cell provides the necessary costimulation signal to T cells to induce a T cell mediated immune response against the transfected tumor cells. In addition, tumor cells which lack MHC class I or MHC class II molecules, or which fail to reexpress sufficient mounts of MHC class I or MHC class II molecules, can be transfected with nucleic acid encoding all or a portion of (e.g., a cytoplasmic-domain truncated portion) of an MHC class I alpha chain protein and beta 2 microglobulin protein or an MHC class II alpha chain protein and an MHC class II beta chain 15 protein to thereby express MHC class I or MHC class II proteins on the cell surface. Expression of the appropriate class I or class II MHC in conjunction with a peptide having the activity of a B lymphocyte antigen (e.g., B7-1, B7-2, B7-3) induces a T cell mediated immune response against the transfected tumor cell. Optionally, a gene encoding an antisense construct which blocks 20 expression of an MHC class II associated protein, such as the invariant chain, can also be cotransfected with a DNA encoding a peptide having the activity of a B lymphocyte antigen to promote presentation of tumor associated antigens and induce tumor specific immunity. Thus, the induction of a T cell mediated immune response in a human subject may be sufficient to overcome tumor-specific tolerance in the subject.

10

25

35

The activity of a protein of the invention may, among other means, be measured by the following methods:

Suitable assays for thymocyte or splenocyte cytotoxicity **30** include, without limitation, those described in: Current Protocols in Immunology, Ed by J. E. Coligan, A. M. Kruisbeek, D. H. Margulies, E. M. Shevach, W. Strober, Pub. Greene Publishing Associates and Wiley-Interscience (Chapter 3, In Vitro assays for Mouse Lymphocyte Function 3.1-3.19: Chapter 7, Immunologic studies in Humans); Herrmann et al., Proc. Natl. Acad. Sci. USA 78:2488-2492, 1981; Herrmann et al., J. Immunol. 128:1968-1974, 1982; Handa et al., J. Immunol. 135:1564-1572, 1985; Takai et al., I. Immunol. 137:3494-3500, 1986; Takai et al., J. Immunol.

140:508-512, 1988; Herrmann et al., Proc. Natl. Acad. Sci. USA 78:2488-2492, 1981; Herrmann et al., J. Immunol. 128:1968-1974, 1982; Handa et al., J. Immunol. 135:1564-1572, 1985; Takai et al., J. Immunol. 137:3494-3500, 1986; Bowmanet al., J. Virology 61:1992-1998; Takai et al., J. Immunol. 140:508-512, 1988; Bertagnolli et al., Cellular Immunology 133:327-341, 1991; Brown et al., J. Immunol. 153:3079-3092, 1994.

Assays for T-cell-dependent immunoglobulin responses and isotype switching (which will identify, among others, proteins that modulate T-cell dependent antibody responses and that affect Thl/Th2 profiles) include, without limitation, those described in: Maliszewski, J. Immunol. 144:3028-3033, 1990; and Assays for B cell function: In vitro antibody production, Mond, J. J. and Brunswick, M. In Current Protocols in Immunology. J. E. e.a. 15 Coligan eds. Vol 1-pp. 3.8.1-3.8.16. John Wiley and Sons. . Toronto. 1994.

10

25

30

35

Mixed lymphocyte reaction (MLR) assays (which will identify, among others, proteins that generate predominantly Thl and CTL responses) include, without limitation, those described in: 20 Current Protocols in Immunology, Ed by J. E. Coligan, A. M. Kruisbeek, D. H. Margulies, E. M. Shevach, W. Strober, Pub. Greene Publishing Associates and Wiley-Interscience (Chapter 3, In Vitro assays for Mouse Lymphocyte Function 3.1-3.19: Chapter 7. Immunologic studies in Humans): Takai et al., J. Immunol. 137:3494-3500, 1986; Takai et al., J. Immunol. 140:508-512, 1988; Bertagnolli et al., J. Immunol. 149:3778-3783, 1992.

Dendritic cell-dependent assays (which will identify, among others, proteins expressed by dendritic cells that activate naive T-cells) include, without limitation, those described in: Guery et al., J. Immunol. 134:536-544, 1995; Inaba et al., Journal of Experimental Medicine 173:549-559, 1991; Macatonia et al., Journal of Immunology 154:5071-5079, 1995; Porgador et al., Journal of Experimental Medicine 182:255-260, 1995; Nair et al., Journal of Virology 67:4062-4069, 1993; Huang et al., Science 264:961-965, 1994; Macatonia et al., Journal of Experimental Medicine 169:1255-1264, 1989; Bhardwaj et al., Journal of Clinical Investigation 94:797-807, 1994; and Inaba et al., Journal of Experimental Medicine 172:631-640, 1990.

Assays for lymphocyte survival/apoptosis (which will identify, among others, proteins that prevent apoptosis after superantigen induction and proteins that regulate lymphocyte homeostasis) include, without limitation, those described in: Darzynkiewicz et al., Cytometry 13:795-808, 1992; Gorczyca et al., Leukemia 7:659-670, 1993; Gorczyca et al., Cancer Research 53:1945-1951, 1993; Itoh et al., Cell 66:233-243, 1991; Zacharchuk, Journal of Immunology 145:4037-4045, 1990; Zamai et al., Cytometry 14:891-897, 1993; Gorczyca et al., International Journal of Oncology 1:639-648, 1992.

10

15

Assays for proteins that influence early steps of T-cell commitment and development include, without limitation, those described in: Antica et al., Blood 84:311-127, 1994; Fine et al., Cellular Immunology 155:111-122, 1994; Galy et al., Blood 85:2770-2778, 1995; Toki et al., Proc. Nat. Acad Sci. USA 88:7548-7551, 1991.

Hematopoiesis Regulating Activity

A protein of the present invention may be useful in regulation of hematopoiesis and, consequently, in the treatment of myeloid or lymphoid cell deficiencies. Even marginal biological activity in support of colony forming cells or of 5 factor-dependent cell lines indicates involvement in regulating hematopoiesis, e.g. in supporting the growth and proliferation of erythroid progenitor cells alone or in combination with other cytokines, thereby indicating utility, for example, in treating various anemias or for use in conjunction with 10 irradiation/chemotherapy to stimulate the production of erythroid precursors and/or erythroid cells; in supporting the growth and proliferation of myeloid cells such as granulocytes and monocytes/macrophages (i.e., traditional CSF activity) useful, for example, in conjunction with chemotherapy to prevent or treat 15 consequent myelo-suppression; in supporting the growth and proliferation of megakaryocytes and consequently of platelets thereby allowing prevention or treatment of various platelet disorders such as thrombocytopenia, and generally for use in place of or complimentary to platelet transfusions; and/or in 20 supporting the growth and proliferation of hematopoietic stem cells which are capable of maturing to any and all of the abovementioned hematopoietic cells and therefore find therapeutic utility in various stem cell disorders (such as those usually treated with transplantation, including, without limitation, 25 aplastic anemia and paroxysmal nocturnal hemoglobinuria), as well as in repopulating the stem cell compartment post irradiation/chemotherapy, either in-vivo or ex-vivo (i.e., in conjunction with bone marrow transplantation or with peripheral progenitor cell transplantation (homologous or heterologous)) as 30 normal cells or genetically manipulated for gene therapy.

The activity of a protein of the invention may, among other means, be measured by the following methods:

Suitable assays for proliferation and differentiation of various hematopoietic lines are cited above.

35

Assays for embryonic stem cell differentiation (which will identify, among others, proteins that influence embryonic differentiation hematopoiesis) include, without limitation, those described in: Johansson et al. Cellular Biology 15:141-151, 1995;

Keller et al., Molecular and Cellular Biology 13:473-486, 1993; McClanahan et al., Blood 81:2903-2915, 1993.

Assays for stem cell survival and differentiation (which will identify, among others, proteins that regulate lymphohematopoiesis) include, without limitation, those described in: Methylcellulose colony forming assays, Freshney, M. G. In Culture of Hematopoietic Cells. R. I. Freshney, et al. eds. Vol pp. 265-268, Wiley-Liss, Inc., New York, N.Y. 1994; Hirayama et al., Proc. Natl. Acad. Sci. USA 89:5907-5911, 1992; Primitive hematopoietic colony forming cells with high proliferative 10 potential, McNiece, I. K. and Briddell, R. A. In Culture of Hematopoietic Cells. R. I. Freshney, et al. eds. Vol pp. 23-39, Wiley-Liss, Inc., New York, N.Y. 1994; Neben et al., Experimental Hematology 22:353-359, 1994; Cobblestone area forming cell assay, Ploemacher, R. E. In Culture of Hematopoietic Cells. R. I. 15 Freshney, et al. eds. Vol pp. 1-21, Wiley-Liss, Inc., New York, N.Y. 1994; Long term bone marrow cultures in the presence of stromal cells, Spooncer, E., Dexter, M. and Allen, T. In Culture of Hematopoietic Cells. R. I. Freshney, et al. eds. Vol pp. 163-20 179, Wiley-Liss, Inc., New York, N.Y. 1994; Long term culture initiating cell assay, Sutherland, H. J. In Culture of Hematopoietic Cells. R. I. Freshney, et al. eds. Vol pp. 139-162, Wiley-Liss, Inc., New York, N.Y. 1994.

Tissue Growth Activity

30

35

A protein of the present invention also may have utility in compositions used for bone, cartilage, tendon, ligament and/or nerve tissue growth or regeneration, as well as for wound healing and tissue repair and replacement, and in the treatment of burns, incisions and ulcers.

A protein of the present invention, which induces cartilage and/or bone growth in circumstances where bone is not normally formed, has application in the healing of bone fractures and cartilage damage or defects in humans and other animals. Such a preparation employing a protein of the invention may have prophylactic use in closed as well as open fracture reduction and also in the improved fixation of artificial joints. De novo bone formation induced by an osteogenic agent contributes to the repair of congenital, trauma induced, or oncologic resection

induced craniofacial defects, and also is useful in cosmetic plastic surgery.

A protein of this invention may also be used in the treatment of periodontal disease, and in other tooth repair processes. Such agents may provide an environment to attract bone-forming cells, stimulate growth of bone-forming cells or induce differentiation of progenitors of bone-forming cells. A protein of the invention may also be useful in the treatment of osteoporosis or osteoarthritis, such as through stimulation of bone and/or cartilage repair or by blocking inflammation or processes of tissue destruction (collagenase activity, osteoclast activity, etc.) mediated by inflammatory processes.

10

15

20

25

30

35

Another category of tissue regeneration activity that may be attributable to the protein of the present invention is tendon/ligament formation--- A protein of the present invention, which induces tendon/ligament-like tissue or other tissue formation in circumstances where such tissue is not normally formed, has application in the healing of tendon or ligament tears, deformities and other tendon or ligament defects in humans and other animals. Such a preparation employing a tendon/ligament-like tissue inducing protein may have prophylactic use in preventing damage to tendon or ligament tissue, as well as use in the improved fixation of tendon or ligament to bone or other tissues, and in repairing defects to tendon or ligament tissue. De novo tendon/ligament-like tissue formation induced by a composition of the present invention contributes to the repair of congenital, trauma induced, or other tendon or ligament defects of other origin, and is also useful in cosmetic plastic surgery for attachment or repair of tendons or ligaments. The compositions of the present invention may provide environment to attract tendon- or ligament-forming cells, stimulate growth of tendon- or ligament-forming cells, induce differentiation of progenitors of tendon- or ligament-forming cells, or induce growth of tendon/ligament cells or progenitors ex vivo for return in vivo to effect tissue repair. The compositions of the invention may also be useful in the treatment of tendonitis, carpal tunnel syndrome and other tendon or ligament defects. The compositions may also include an

PCT/IB01/02050 WO 01/98454

appropriate matrix and/or sequestering agent as a carrier as is well known in the art-

The protein of the present invention may also be useful for proliferation of neural cells and for regeneration of nerve and 5 brain tissue, i.e. for the treatment of central and peripheral nervous system diseases and neuropathies, as well as mechanical and traumatic disorders, which involve degeneration, death or trauma to neural cells or nerve tissue. More specifically, a protein may be used in the treatment of diseases of the peripheral nervous system, such as peripheral nerve injuries, peripheral neuropathy and localized neuropathies, and central nervous system diseases, such as Alzheimer's, Parkinson's disease, Huntington's disease, amyotrophic lateral sclerosis, and Shy-Drager syndrome. Further conditions which may be treated in 15 accordance—with the present invention include mechanical and traumatic disorders, such as spinal cord disorders, head trauma and cerebrovascular diseases such as stroke. Peripheral neuropathies resulting from chemotherapy or other medical therapies may also be treatable using a protein of the invention.

10

20

25

35

Proteins of the invention may also be useful to promote better or faster closure of non-healing wounds, including without. limitation pressure ulcers, ulcers associated with vascular insufficiency, surgical and traumatic wounds, and the like.

It is expected that a protein of the present invention may also exhibit activity for generation or regeneration of other tissues, such as organs (including, for example, pancreas, liver, intestine, kidney, skin, endothelium), muscle (smooth, skeletal or cardiac) and vascular (including vascular endothelium) tissue, or for promoting the growth of cells comprising such tissues. 30 Part of the desired effects may be by inhibition or modulation of fibrotic scarring to allow normal tissue to regenerate. A protein of the invention may also exhibit angiogenic activity.

A protein of the present invention may also be useful for gut protection or regeneration and treatment of lung or liver fibrosis, reperfusion injury in various tissues, and conditions resulting from systemic cytokine damage.

A protein of the present invention may also be useful for promoting or inhibiting differentiation of tissues described

PCT/IB01/02050 WO 01/98454

above from precursor tissues or cells; or for inhibiting the growth of tissues described above.

The activity of a protein of the invention may, among other means, be measured by the following methods:

Assays for tissue generation activity include, without limitation, those described in: International Patent Publication No. W095/16035 (bone, cartilage, tendon); International Patent Publication No. W095/05846 (nerve, neuronal); International - Patent Publication No. W091/07491 (skin, endothelium).

Assays for wound healing activity include, without limitation, those described in: Winter, Epidermal Wound Healing, pps. 71-112 (Maibach, H. I. and Rovee, D. T., eds.), Year Book Medical Publishers, Inc., Chicago, as modified by Eaglstein and Mertz, J. Invest. Dermatol 71:382-84 (1978).

15 Activin/Inhibin Activity

5

10

25

30

35

A protein of the present invention may also exhibit activinor inhibin-related activities. Inhibins are characterized by their ability to inhibit the release of follicle stimulating hormone (FSH), while activins and are characterized by their 20 ability to stimulate the release of follicle stimulating hormone (FSH). Thus, a protein of the present invention, alone or in heterodimers with a member of the inhibin alpha family, may be useful as a contraceptive based on the ability of inhibins to decrease fertility in female mammals and decrease spermatogenesis in male mammals. Administration of sufficient amounts of other inhibins can induce infertility in these mammals. Alternatively, the protein of the invention, as a homodimer or as a heterodimer with other protein subunits of the inhibin- beta group, may be useful as a fertility inducing therapeutic, based upon the ability of activin molecules in stimulating FSH release from cells of the anterior pituitary. See, for example, U.S. Pat. No. 4.798.885. A protein of the invention may also be useful for advancement of the onset of fertility in sexually immature mammals, so as to increase the lifetime reproductive performance of domestic animals such as cows, sheep and pigs.

The activity of a protein of the invention may, among other means, be measured by the following methods:

Assays for activin/inhibin activity include, without limitation, those described in: Vale et al., Endocrinology 91:562-572, 1972; Ling et al., Nature 321:779-782, 1986; Vale et al., Nature 321:776-779, 1986; Mason et al., Nature 318:659-663, 1985; Forage et al., Proc. Natl. Acad. Sci. USA 83:3091-3095, 1986.

Chemotactic/Chemokinetic Activity

10

15

20

25

30

35

A protein of the present invention may have chemotactic or chemokinetic activity (e.g., act as a chemokine) for mammalian cells, including, for example, monocytes, fibroblasts, neutrophils, T-cells, mast cells, eosinophils, epithelial and/or endothelial cells. Chemotactic and chemokinetic proteins can be used to mobilize or attract a desired cell population to a desired site of action. Chemotactic or chemokinetic proteins provide particular advantages in treatment of wounds and other trauma to tissues, as well as in treatment of localized infections. For example, attraction of lymphocytes, monocytes or neutrophils to tumors or sites of infection may result in improved immune responses against the tumor or infecting agent.

A protein or peptide has chemotactic activity for a particular cell population if it can stimulate, directly or indirectly, the directed orientation or movement of such cell population. Preferably, the protein or peptide has the ability to directly stimulate directed movement of cells. Whether a particular protein has chemotactic activity for a population of cells can be readily determined by employing such protein or peptide in any known assay for cell chemotaxis.

The activity of a protein of the invention may, among other means, be measured by the following methods:

Assays for chemotactic activity (which will identify proteins that induce or prevent chemotaxis) consist of assays that measure the ability of a protein to induce the migration of cells across a membrane as well as the ability of a protein to induce the adhesion of one cell population to another cell population. Suitable assays for movement and adhesion include, without limitation, those described in: Current Protocols in Immunology, Ed by J. E. Coligan, A. M. Kruisbeek, D. H. Marguiles, E. M. Shevach, W. Strober, Pub. Greene Publishing

Associates and Wiley-Interscience (Chapter 6.12, Measurement of alpha and beta Chemokines 6.12.1-6.12.28; Taub et al. J. Clin. Invest. 95:1370-1376, 1995; Lind et al. APMIS 103:140-146, 1995; Muller et al Eur. J. Immunol. 25:1744-1748; Gruber et al. J. of Immunol. 152:5860-5867, 1994; Johnston et al. J. of Immunol. 153:1762-1768, 1994.

Hemostatic and Thrombolytic Activity

A protein of the invention may also exhibit hemostatic or thrombolytic activity. As a result, such a protein is expected to 10 · be useful in treatment of various coagulation disorders (including hereditary disorders, such as hemophilias) or to enhance coagulation and other hemostatic events in treating wounds resulting from trauma, surgery or other causes. A protein of the invention may also be useful for dissolving or inhibiting formation of thromboses and for treatment and prevention of conditions resulting therefrom (such as, for example, infarction of cardiac and central nervous system vessels (e.g., stroke).

The activity of a protein of the invention may, among other means, be measured by the following methods:

Assay for hemostatic and thrombolytic activity include, without limitation, those described in: Linet et al., J. Clin. Pharmacol. 26:131-140, 1986; Burdick et al., Thrombosis Res. .45:413-419.1987:_Humphrey et al., Fibrinolysis 5:71-79 (1991); Schaub, Prostaglandins 35:467-474, 1988.

25 Receptor/Ligand Activity

15

20

30

35

A protein of the present invention may also demonstrate activity as receptors, receptor ligands or inhibitors or agonists of receptor/ligand interactions. Examples of such receptors and ligands include, without limitation, cytokine receptors and their ligands, receptor kinases and their ligands, receptor phosphatases and their ligands, receptors involved in cell-cell interactions and their ligands (including without limitation, cellular adhesion molecules (such as selectins, integrins and their ligands) and receptor/ligand pairs involved in antigen presentation, antigen recognition and development of cellular and humoral immune responses). Receptors and ligands are also useful for screening of potential peptide or small molecule inhibitors

of the relevant receptor/ligand interaction. A protein of the present invention (including, without limitation, fragments of receptors and ligands) may themselves be useful as inhibitors of receptor/ligand interactions.

The activity of a protein of the invention may, among other means, be measured by the following methods:

Suitable assays for receptor-ligand activity include without limitation those described in:Current Protocols in Immunology, Ed by J. E. Coligan, A. M. Kruisbeek, D. H. Margulies, E. M. Shevach, W. Strober, Pub. Greene Publishing Associates and Wiley-Interscience (Chapter 7.28, Measurement of Cellular Adhesion under static conditions 7.28.1-7.28.22), Takai et al., Proc. Natl. Acad. Sci. USA 84:6864-6868, 1987; Bierer et al., J. Exp. Med. 168:1145-1156, 1988; Rosenstein et al., J. Exp. Med. 169:149-160 1989; Stoltenborg et al., J. Immunol. Methods 175:59-68, 1994; Stitt et al., Cell 80:661-670, 1995.

Anti-Inflammatory Activity

10

15

Proteins of the present invention may also exhibit antiinflammatory activity. The anti-inflammatory activity may be 20 achieved by providing a stimulus to cells involved in the inflammatory response, by inhibiting or promoting cell-cell interactions (such as, for example, cell adhesion), by inhibiting or promoting chemotaxis of cells involved in the inflammatory process, inhibiting or promoting cell extravasation, or by 25 stimulating or suppressing production of other factors which more directly inhibit or promote an inflammatory response. Proteins exhibiting such activities can be used to treat inflammatory conditions including chronic or acute conditions), including without limitation intimation associated with infection (such as 30 septic shock, sepsis or systemic inflammatory response syndrome (SIRS)), ischemia-reperfusion injury, endotoxin lethality, arthritis, complement-mediated hyperacute rejection, nephritis, cytokine or chemokine-induced lung injury, inflammatory bowel disease. (rohn's disease or resulting from over production of cytokines such as TNF or IL-1. Proteins of the invention may also 35 be useful to treat anaphylaxis and hypersensitivity to an antigenic substance or material.

Tumor Inhibition Activity

In addition to the activities described above for immunological treatment or prevention of tumors, a protein of the invention may exhibit other anti-tumor activities. A protein may inhibit tumor growth directly or indirectly (such as, for example, via ADCC). A protein may exhibit its tumor inhibitory activity by acting on tumor tissue or tumor precursor tissue, by inhibiting formation of tissues necessary to support tumor growth (such as, for example, by inhibiting angiogenesis), by causing production of other factors, agents or cell types which inhibit tumor growth, or by suppressing, eliminating or inhibiting factors, agents or cell types which promote tumor growth.

Other Activities

10

A protein of the invention may also exhibit one or more of the following additional activities or effects: inhibiting the 15 growth, infection or function of, or killing, infectious agents, including, without limitation, bacteria, viruses, fungi and other parasites; effecting (suppressing or enhancing) bodily characteristics, including, without limitation, height, weight, 20 hair color, eye color, skin, fat to lean ratio or other tissue pigmentation, or organ or body part size or shape (such as, for example, breast augmentation or diminution, change in bone form or shape); effecting biorhythms or caricadic cycles or rhythms; effecting the fertility of male or female subjects; effecting the 25 metabolism, catabolism, anabolism, processing, utilization, storage or elimination of dietary fat, lipid, protein, carbohydrate, vitamins, minerals, cofactors or other nutritional factors or component(s); effecting behavioral characteristics, including, without limitation, appetite, libido, stress, 30 cognition (including cognitive disorders), depression (including depressive disorders) and violent behaviors; providing analgesic effects or other pain reducing effects; promoting differentiation and growth of embryonic stem cells in lineages other than hematopoietic lineages; hormonal or endocrine activity; in the 35 case of enzymes, correcting deficiencies of the enzyme and treating deficiency-related diseases; treatment of hyperproliferative disorders (such as for example, psoriasis); immunoglobulin-like activity (such as, for example, the ability

to bind antigens or complement); and the ability to act as an antigen in a vaccine composition to raise an immune response against such protein or another material or entity which is cross-reactive with such protein.

5 Particular Applications for Certain Clones

The following sets out a non-exclusive list of applications for certain embodiments of the invention. In the interest of economy, applications relevant to multiple embodiments are not duplicated in this list. Other embodiments described herein have similar characteristics, as described there. The artisan is directed, therefore, to the Description of the Sequences for similar descriptions of the functions of other embodiment.

Testes

10

- htes3_lOilb: The new protein can find application in diagnosis/therapy in leukemia predisposition/disease in the modulation of DNA repair.
- htes3_10n10: The new protein can find application in studying the expression profile of testis-specific genes.

htes3_llal7: The new protein can find application in studying_the_expression profile_of testis-specific genes and as a new marker for testicular cells-

- 25
 htes3_llc22: The new protein can find application in
 modulating/blocking of regulatory pathways.
- htes3_lld2l: The new protein can find application in diagnosis of diseases due to unnormal protein degradation like muscular dystrophy or multiple sclerosis as well as in modulating the half life of specific proteins and in expression profiling.

35 Kidney

hfkd2_3k1 The new protein can find application in modulation of endocytosis.strong similarity to testicular dynamin (Rattus norvegicus).

Amygdala:

5

10

15

35

40

hamy2_10h17: The new protein can find application in modulating protein-protein-interaction and in studying the expression profile of amygdala-specific genes.

hamy2_10p7: The new protein can find application in modulation of NA+/Ca2+-exchange and voltage-dependend processes.

hamy2_lld2: The new protein can find application in studying the expression profile of amygdala-specific genes and as a new marker for amygdala cells.

hamy2_lln4: The new protein can find application in modulation of DNA-repair and a as a new tool for manipulation of nucleic acids.

20 hamy2_121f19: The new protein can find application modulation of cyto skeleton-membrane interactions.

Fetal Brain:

hfbr2_78cl2: The new protein can find application in the modulation of translational pathways.

hfbr2_78dl8: The new protein can find application in studying the expression profile of brain-specific genes.

30 hfbr2_78d4: The new protein can find application in studying the expression profile of brain-specific genes and as a new marker for amygdala cells.

hfbr2_78el8: The new protein can find application in studying the expression profile of brain-specific genes.

hfbr2_78i21: The new protein can find application in diagnosis/modulation of protein damage and age-related degenerative processes.

Melanoma:

hmel2_l2jl: The new protein can find application in studying the expression profile of melanoma-specific genes.

hmel2_7g14: The new protein can find application in modulation of the sorting of proteins into different compartments.

hmel2_7kl9: The new protein can find application in studying the expression profile of melanoma-specific genes.

10

15

20

25

30

35

5

VARIANTS OF THE INVENTIVE DNA MOLECULES

Variants in General

"Variants," according to the invention, include DNA and/or protein molecules that resemble, structurally and/or functionally, those set forth herein. Variants may be isolated from natural sources ("homologs"), may be entirely synthetic or may be based in part on both natural and synthetic approaches.

The section set forth below presents various structural and functional characteristics of molecules within the invention. Preferred molecules are characterized by a combination of one or more of these characteristics. For instance, some preferred molecules-are-described-with reference to-at least two structural characteristics, while others may be described with reference to structural least one and at least one functional characteristic.

It will be recognized by the skilled artisan that structure ultimately defines function, i.e. the functions of the molecules described herein derives from the structures of those molecules. Accordingly, the structural variants described below that bear the closest structural relationship (as variously defined below) to the inventive molecules are the variants that most likely will preserve biological function. This relationship between structure and function will guide the skilled artisan in identifying the preferred embodiments of the invention.

Splicing Variants

10

15

20

25

30

35

It is well-known that eukaryotic structural genes are comprised of both protein coding and non-coding portions. When the messenger RNA is transcribed from the DNA template it contains introns, which are non-coding, and exons, which are coding. In order to form a translation competent mRNA, the introns must be "spliced" out of this initial pre mRNA.

Specific sequences within the pre mRNA represent "splice junctions" that direct the cellular splicing machinery to the appropriate position. The splice junctions are loosely conserved sequence regions of the pre mRNA, which almost invariably begin with GT and end with AG (DNA perspective). The 5' end of the splice junction typically contains about nine somewhat conserved residues, for example, C/AAGTA/GAGT. The 3' end usually contains a pyrimidine rich stretch of at least about 11 nucleotides, followed by NC/TAGG. Splicing occurs before the GT and after the AG. Mount, Nucleic Acids Res. 10:459-72 (1982).

Interestingly, exons often correspond to discrete functional domains of the protein product. The intron/exon arrangement thus creates a linear array of nucleotides which can be correlated to discrete, and often interchangeable, functional protein fragments. Go. Nature 291:90-92 (1981); Branden et al., EMBO J. 3:1307-10 (1984). This linear arrangement creates the possibility of generating multiple different full length proteins by rearranging the order of the different functional portions in the array. For example, if a set of exons are arranged 1-2-3-4, where (-) represents the introns separating the exons, a splicing event need not simply produce 1234, but may produce 123, 134, 124 and so on. Production of different mRNA products in this way is commonly called "alternative splicing." Andreadis et al., Ann. Rev. Cell Biol. 3:207-42 (1987).

Some of the present DNA molecules can be represented in modular fashion in terms of their coding regions. Essentially, these modules are exons (though each "exon" may in fact be made up of several exons), which may be combined in different ways to form a variety of different DNA molecules, each encoding a different functional protein. Splicing variants are indicated in the Description of the Sequences.

Degenerate Variants

10

---15

20

25

30

35

One aspect of the present invention provides "degenerate variants" of the nucleic acid fragments of the present invention. A "degenerate variant" is a nucleotide fragment which differs from those of inventive molecules by nucleotide sequence, but due to the degeneracy of the genetic code, encodes an identical polypeptide sequence.

Given the known relationship between DNA sequences and the proteins they encode, degenerate variants typically are described by reference to this relationship. It is well known that the degeneracy of the genetic code results in many possible DNA sequences which encode a particular protein. Indeed, of the three bases which comprise an amino acid-encoding triplet, the third position, and often the second, almost always may vary. This fact alone allows for a class of variant DNA molecules which encode protein sequences identical to those disclosed herein, yet have about 30% sequence variation. In other words, the variant DNA molecules are about 70% identical to the inventive DNAs, having no additional or deleted sequences. Thus, one aspect invention provides degenerate variant DNA molecules encoding the inventive protein sequences.

In one embodiment, these variants have at least about 70% sequence identity with the DNA molecules described herein. In a preferred embodiment, these variants have at least about 80% sequence identity to the inventive molecules. In a more preferred embodiment these variants have at least about 90% sequence identity with the inventive molecules.

Conservative Amino Acid Variants

Variants according to the invention also may be made that conserve the overall molecular structure of the encoded proteins. Given the properties of the individual amino acids comprising the disclosed protein products, some rational substitutions will be recognized by the skilled worker. Amino acid substitutions, i.e. "conservative substitutions," may be made, for instance, on the basis of similarity in polarity, charge, solubility, hydrophobicity, hydrophilicity, and/or the amphipathic nature of the residues involved.

For example: (a) nonpolar (hydrophobic) amino acids include alanine, leucine, isoleucine, valine, proline, phenylalanine, tryptophan, and methionine; (b) polar neutral amino acids include glycine, serine, threonine, cysteine, tyrosine, asparagine, and glutamine; (c) positively charged (basic) amino acids include arginine, lysine, and histidine; and (d) negatively charged (acidic) amino acids include aspartic acid and glutamic acid. Substitutions typically may be made within groups (a)-(d). addition, glycine and proline may be substituted for one another based on their ability to disrupt α -helices. Similarly, certain amino acids, such as alanine, cysteine, leucine, methionine, glutamic acid, glutamine, histidine and lysine are more commonly found in α -helices, while valine, isoleucine, phenylalanine, tyrosine, tryptophan and threonine are more commonly found in β pleated sheets. Glycine, serine, aspartic acid, asparagine, and proline are commonly found in turns. Some preferred substitutions may be made among the following groups: (i) S and T; (ii) P and G; and (iii) A₁ V₁ L and I. Given the known genetic code, and recombinant and synthetic DNA techniques, the skilled scientist readily can construct DNAs encoding the conservative amino acid variants.

As used herein, "sequence identity" between two polypeptide sequences indicates the percentage of amino acids that are identical between the sequences. "Sequence similarity" indicates the percentage of amino acids that either are identical or that represent conservative amino acid substitutions.

Functionally Equivalent Variants

10

15

20

25

30

35

Yet another class of DNA variants within the scope of the invention may be described with reference to the product they encode. As shown in the Description of the Sequences, some of the inventive DNA molecules encode a protein having a degree of homology with known proteins, or protein domains. It is expected, therefore, that they will have some or all of the requisite functional features of such molecules. These "functionally equivalent variants" products are characterized by the fact that they are functionally equivalent, with respect to biological activity, to certain known molecules.

Also provided herein is information on common structural motifs, including consensus sequences that will guide the artisan in constructing functionally equivalent variants. It will be understood that the motifs, identified in the Description of the Sequences for each inventive protein, may be modified within the identified consensus sequences. Thus, the invention contemplates the proteins in the Description of the Sequences that contain variability in the consensus sequences identified, and the invention further contemplates the full range of nucleic acids encoding them, and the complements of those nucleic acids.

Hybridizing Variants

10

15

20

25

30

35

DNA variants within the invention also may be described by reference to their physical properties in hybridization. skilled in the field will recognize that DNA can be used to identify its complement and, since DNA is double stranded, its or homologa using nucleic acid hybridization techniques. It will also be recognized that hybridization can occur with less than 100% complementarity. However appropriate choice of conditions, hybridization techniques can be used to differentiate among DNA sequences based structural relatedness to a particular proberegarding such conditions see, for example, Sambrook et al., 1989, MOLECULAR CLONING, A LABORATORY MANUAL, Cold Spring Harbor Press, N.Y.; and Ausubel et al., 1989, CURRENT PROTOCOLS IN MOLECULAR BIOLOGY, Green Publishing Associates and Wiley Interscience, N.Y.

Structural relatedness between two polynucleotide sequences can be expressed as a function of "stringency" of the conditions under which the two sequences will hybridize with one another. As used herein, the term "stringency" refers to the extent that the conditions disfavor hybridization. Stringent conditions strongly disfavor hybridization, and only the most structurally related molecules will hybridize to one another under such conditions. Conversely, non-stringent conditions favor hybridization of molecules displaying a lesser degree of structural relatedness. Hybridization stringency, therefore, directly correlates with the structural relationships of two nucleic acid sequences. The following relationships are useful in correlating hybridization

-555-

and relatedness (where T_m is the melting temperature of a nucleic acid duplex):

a. $T_m = 69.3 + 0.41(6+0)$ %

5

15

35

40

- b. The T_m of a duplex DNA decreases by $1^{\circ}C$ with every increase of 1% in the number of mismatched base pairs.
- c. $(T_m)_{\mu^2}$ $(T_m)_{\mu^3}$ = 18.5 $\log_{10}\mu^2/\mu^3$ where μ^3 and μ^2 are the ionic strengths of two solutions.

Hybridization stringency is a function of many factors, including overall DNA concentration, ionic strength, temperature, probe size and the presence of agents which disrupt hydrogen bonding. Factors promoting hybridization include high DNA concentrations, high ionic strengths, low temperatures, longer probe size and the absence of agents that disrupt hydrogen bonding.

Hybridization usually is done in two stages. 20 First, in the "binding" stage, the probe is bound to the target under conditions favoring hybridization. Stringency is usually controlled at this stage by altering the temperature. For high stringency, the temperature is usually between 65°C and 70°C, unless short (<20 25 nt) oligonucleotide probes are used. Α representative hybridization solution comprises LX SSC, D.5% SDS, 5X Denhardt's solution and 100μg of non-specific carrier DNA. See Ausubel et al., supra, section 2.9, supplement 27 (1994). Of course many different, yet functionally equivalent, buffer conditions 30 known. Where the degree of relatedness is lower, a lower temperature may be chosen. Low stringency binding temperatures are between about 25°C and 40°C. Medium stringency is between at least about 40°C to less than about 65°C. High stringency is at least about 65°C.

Second, the excess probe is removed by washing. It is at this stage that more stringent conditions usually are applied. Hence, it is this "washing" stage that is most important in determining relatedness via hybridization. Washing solutions typically contain lower salt concentrations. One exemplary medium stringency solution contains 2X SSC and 0.1% SDS. A high stringency wash solution contains the equivalent (in ionic

strength) of less than about 0.2X SSC, with a preferred stringent solution containing about 0.1X SSC. The temperatures associated with various stringencies are the same as discussed above for "binding." The washing solution also typically is replaced a number of times during washing. For example, typical high stringency washing conditions comprise washing twice for 30 minutes at 55° C. and three times for 15 minutes at 60° C.

The present invention includes nucleic acid molecules that hybridize to the inventive molecules under high stringency binding and washing conditions. More preferred molecules (from an mRNA perspective) are those that are at least 50 % of the length of any one of those depicted in the Description of the Sequences. Particularly preferred molecules are at least 75 % of the length of those molecules.

15 Substitutions, Insertions, Additions and Deletions

10

20

25

30

35

In a general sense, the preferred DNA variants of the invention are those that retain the closest relationship, as described by "sequence identity" to the inventive DNA molecules. According to another aspect of the invention, therefore, substitutions, insertions, additions and deletions of defined properties are contemplated. It will be recognized that sequence . identity between two polynucleotide sequences, as defined herein, generally is determined with reference to the protein coding region of the sequences. Thus, this definition does not at all limit the amount of DNA, such as vector DNA, that may be attached to the molecules described herein. Preferred DNA sequence variants include molecules encoding proteins sharing some or all of any relevant biological activity of the native molecule.

In creating these variants, the skilled worker will be guided by reference to the protein structure. First, insertions and deletions in any recognized functional domain above generally should be avoided, except as noted below in the section entitled "Proteins," where this domain is discussed in detail. Alterations in such domains usually will be limited to conservative amino acid substitutions. In addition, where insertions and deletions are desired, this may be accomplished at the N- and/or C-terminus of the protein molecule (or the corresponding coding regions of the DNA). If insertions or deletions are made within the protein.

-557-

WO 01/98454

20

25

30

35

PCT/IB01/02050

deletions of major structural features usually should be avoided. Thus, a preferred place to make insertion or deletion variants is in non-structural regions, such as linker regions between two alpha helices.

"Substitutions" generally refer to alterations in the DNA 5 sequence which do not change its overall length, but only alter one or more nucleotide positions, substituting one for another in the common sense of the word. One class οf substitutions, "degenerate substitutions," are those that do not alter the encoded amino acid sequence. Some substitutions retains 10 50%, 55%, 60% or 65% identity. Preferred substitutions retain at least about 70% identity, more preferably at least 70% or 75% identity, with the inventive DNAs. Some more preferred molecules have at least about 80% identity, more preferably at least 80% or 15 - 85% identity. Particularly preferred DNAs share at least about 90% identity, more preferably at least 90% or 95% identity.

"Insertions," unlike substitutions, alter the overall length of the DNA molecule, and thus sometimes the encoded protein. Insertions add extra nucleotides to the interior (not the 5' or 3' ends) of the subject DNAs. Preferred insertions are made with reference to the protein sequence encoded by the DNA. Thus, it is most preferred to provide an insertion in the DNA at a location that corresponds to an area of the encoded protein which lacks structure. For instance, it typically would not be beneficial, if the preservation of biological activity is desired, to provide an insertion within an alpha-helical region or a beta-pleated sheet. Accordingly, non-structural areas, such as those containing helix-breaking glycines and proline residues, are most preferred sites of insertion. Other preferred sites of insertion are the splice sites, which are indicated above in the description of the inventive DNA molecules.

While the optimal size of insertions will vary depending upon the site of insertion and its effect on the overall conformation of the encoded protein, some general guides are useful. Generally, the total insertions (irrespective of their number) should not add more than about 30% (or preferably not more than 30%) to the overall size of the encoded protein. More preferably, the insertion adds less than about 10~20% (yet more preferably 10~20%) in size, with less than about 10% being most preferred. The

number of insertions is limited only by the number of suitable insertions sites, and secondarily by the foregoing preferences.

"Additions," like insertions, also add to the overall size of the DNA molecule, and usually the encoded protein. instead of being made within the molecule, they are made on the 5° or 3' end, usually corresponding to the N- or C- terminus of the encoded protein. Unlike deletions, additions are not very size-Indeed, additions may be of virtually any size. Preferred additions, however, do not exceed about 100% of the size of the native molecule. More preferably, they add less than about 60 to 30% to the overall size, with less than about 30% being most preferred.

10

15

20

25

30

35

"Deletions" diminish the overall size of the DNA and therefore, also reduce the size of the protein encoded by that Deletions may be made from either end of the molecule or Typical preferred deletions remove discrete internal to itstructural features of the encoded protein. For example, some deletions will comprise the deletion of one or more exons which may define a structural feature. Preferred deletions remove less than about 30% of the size of the subject molecule. More preferred deletions remove less than about 20% and most preferred deletions remove less than about 10%.

----Computer-Defined Variants and Definition of "Sequence Identity"

In general, both the DNA and protein molecules of the invention can be defined with reference to "sequence identity." As used herein, "sequence identity" refers to a comparison made between two molecules using, for example, the standard Smith-Waterman algorithm that is well known in the art.

Some molecules have at lease about 50%, 55% or 60% identity. Preferred molecules are those having at least about 65% sequence identity, more preferably at least 65% or 70% sequence identity. Other preferred molecules have at least about 80%, more preferably at least 80% or 85%, sequence identity. Particularly preferred molecules have at least about 90% sequence identity, more preferably at least 90% sequence identity. Most preferred molecules have at least about 95%, more preferably at least 95%, sequence identity. As used herein, two nucleic acid molecules or

proteins are said to "share significant sequence identity" if the two contain regions which possess greater than 85% sequence (amino acid or nucleic acid) identity.

"Sequence identity" is defined herein with reference the Blast 2 algorithm, which is available at the NCBI (http://www.ncbi.nlm.nih.gov/BLAST), using default parameters. References pertaining to this algorithm include: those found at http://www.ncbi.nlm.nih.gov/BLAST/blast_references.html; Altschul, S.F., Gish, W., Miller, W., Myers, E.W. & Lipman, D.J. (1990) "Basic local alignment search tool." J. Mol. Biol. 10 215:403-410; Gish, W. & States, D.J. (1993) "Identification of protein coding regions by database similarity search." Nature Genet. 3:266-272; Madden, T.L., Tatusov, R.L. & Zhang, J. (1996) "Applications of network BLAST server" Meth. Enzymol. 266:131-15 141; Altschul, S.F., Madden, T.L., Schäffer, A.A., Zhang, J., Zhang, Z., Miller, W. & Lipman, D.J. (1997) "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs." Nucleic Acids Res. 25:3389-3402; and Zhang, J. & Madden, T.L. (1997) "PowerBLAST: A new network BLAST application for 20 interactive or automated sequence analysis and annotation." Genome Res. 7:649-656.

METHODS OF MAKING VARIANTS

25

30

It will be recognized that variants of the inventive molecules can be constructed in several different ways. For example, they may be constructed as completely synthetic DNAs. Methods of efficiently synthesizing oligonucleotides in the range of 20 to about 150 nucleotides are widely available. See Ausubel et al., supra, section 2.11, Supplement 21 (1993). Overlapping oligonucleotides may be synthesized and assembled in a fashion first reported by Khorana et al., J. Mol. Biol. 72:209-217 (1971); see also Ausubel et al. Section 8.2. The synthetic DNAs are designed with convenient restriction sites engineered at the 5' and 3' ends of the gene to facilitate cloning into an appropriate vector.

An alternative method of generating variants is to start with one of the inventive DNAs and then to conduct site-directed mutagenesis. See Ausubel et al., supra, chapter & Supplement 37

(1997). In a typical method, a target DNA is cloned into a single-stranded DNA bacteriophage vehicle. Single-stranded DNA is isolated and hybridized with a oligonucleotide containing the desired nucleotide alteration(s). The complementary strand is synthesized and the double stranded phage is introduced into a host. Some of the resulting progeny will contain the desired mutant, which can be confirmed using DNA sequencing. In addition, various methods are available that increase the probability that the progeny phage will be the desired mutant. These methods are well known to those in the field and kits are commercially available for generating such mutants.

ISOLATING HOMOLOGS

Methods

10

15

20

25

30

By using the sequences disclosed herein as probes or as primers, and techniques such as PCR cloning and colony/plaque hybridization, one skilled in the art can obtain homologs. "Homologs" are essentially naturally-occurring variants and include allelic, species-specific and tissue-specific variants.

Region-specific primers or probes derived from the nucleotide sequence(s) provided can be used to prime DNA synthesis and PCR amplification, as well as to identify colonies containing cloned DNA encoding a homolog using known methods (Innis et al., PCR Protocols, Academic Press, San Diego, CA (1990)). Such an application is useful in diagnostic methods, as described in more detail below, as well as in preparing full-length DNAs from various sources. The PCR primers are preferably at least 15 bases, and more preferably at least 16 bases in length. When selecting a primer sequence, it is preferred that the primer pairs have approximately the same G/C ratio, so that melting temperatures are approximately the same. As a general guide, the formula $3(G+C) + 2(A+T) = {}^{\circ}C_{1}$ is useful.

When using primers derived from the inventive sequences, one skilled in the art will recognize that by employing high stringency conditions (e.g., annealing at 50-60°C), only sequences with greater than 75% sequence identity to the primer will be amplified. By employing lower stringency conditions (e.g.,

-561-

annealing at 35-37°C), sequences which have greater than 40-50% sequence identity to the primer also will be amplified.

The PCR product may be subcloned and sequenced to confirm that it indeed displays the expected sequence identity. fragment may then be used to isolate a full length cDNA clone by a variety of methods. For example, the amplified fragment may be labeled and used to screen a bacteriophage cDNA Alternatively, the labeled fragment may be used to screen a genomic library.

10

20

25

30

35

PCR technology may also be utilized to isolate full length For example, RNA may be isolated, following standard procedures, from an appropriate cellular or tissue source. A reverse transcription reaction may be performed on the RNA using an oligonucleotide primer specific for the most 5' end 15 of the amplified fragment for the priming of first strand synthesis. The resulting RNA/DNA hybrid may then be "tailed" with guanines using a standard terminal transferase reaction, the hybrid may be digested with RNAase H₁ and second strand synthesis may then be primed with a poly-C primer. Thus, cDNA sequences upstream of the amplified fragment may easily be isolated. For a review of cloning strategies which may be used, see e.g., Sambrook et al., 1989, supra.

When using DNA probes derived from the inventive sequences for colony/plaque hybridization, one skilled in the art will recognize that by employing medium to high stringency conditions (e.g., hybridizing at 50-65°C in 5X SSPC and 50% formamide, and washing at 50-65°C in 0.5% SSPC), sequences having regions with greater than 90% sequence identity to the probe can be obtained. that by employing lower stringency conditions hybridizing at 35-37°C in 5X SSPC and 40-45% formamide, washing at 42°C in SSPC), sequences having regions with greater than 35-45% sequence identity to the probe will be obtained.

Suitably, genomic or cDNA libraries can be constructed and screened in accord with the previous paragraph. The libraries should be derived from a tissue or organism that is known to express the gene of interest, or that is suspected of expressing the gene. The clone containing the homolog may then be purified

through methods routinely practiced in the art, and subjected to sequence analysis.

Additionally, an expression library can be constructed utilizing DNA isolated from or cDNA synthesized from a tissue or organism that is known to express the gene of interest, or that is suspected of expressing the gene. In this manner, clones may be induced and screened using standard antibody screening techniques in conjunction with antibodies raised against the normal gene product, as described herein. (For screening techniques, see, for example, Harlow, E. and Lane, eds., 1988, ANTIBODIES: A LABORATORY MANUAL, Cold Spring Harbor Press.)

Human Homologs

10

15

25

30

35

Any organism or tissue can be used as the source for homologs of the present invention so long as the organism or tissue naturally expresses such a protein or contains genes encoding the same. The most preferred organism for isolating homologs is human.

PROTEINS OF THE INVENTION

20 One class of proteins included within the invention is encoded by the inventive DNA molecules presented. Other proteins according to the invention are those encoded by the DNA variants described above. As noted, these variants are designed with the encoded proteins in mind.

A preferred class of protein fragments includes those fragments which retain any biological activity. These molecules share functional features common the family of proteins, although these characteristics may vary in degree.

According to one aspect of the invention fragments of the inventive proteins are contemplated. Some preferred fragments are those which are capable of eliciting an immune response. Generally these "antigenic" fragments will be from about five amino acids in length to about fifty amino acids in length. Some preferred antigenic fragments are from five to about twenty amino acids long. "Antigenic" response may refer to a T cell response a B cell response or a response by cells of the macrophage/monocyte lineages. In most cases, however, it will

refer to the immune response involved in the generation of antibodies. In other words, the relevant immune response is that of helper T cells and/or B cells. These preferred molecules comprise one or more T cell and /or B cell epitopes.

5 ANTIBODIES OF THE INVENTION

10

15

20

25

30

35

Antibodies raised against the proteins and protein fragments of the invention also are contemplated by the invention. Described below are antibody products and methods for producing antibodies capable of specifically recognizing one or more epitopes of the presently described proteins and their derivatives.

Antibodies include but are not limited to polyclonal antibodies monoclonal antibodies (mAbs) humanized or chimeric antibodies single chain antibodies including single chain Fv (scFv) fragments. Fab fragments F(ab') fragments fragments produced by a Fab expression library anti-idiotypic (anti-Id) antibodies epitope-binding fragments and humanized forms of any of the above.

As known to one in the art, these antibodies may be used, for example, in the detection of a target protein in a biological sample. They also may be utilized as part of treatment methods, and/or may be used as part of diagnostic techniques whereby patients may be tested for abnormal levels or for the presence of abnormal forms of the such proteins.

general, techniques for preparing polyclonal monoclonal antibodies as well as hybridomas capable of producing the desired antibody are well known in the art (Campbell, A.M., Antibody Technology: Laboratory Techniques Biochemistry and Molecular Biology, Elsevier Science Publishers, Amsterdam, The Netherlands (1984); St. Groth et al., J. Immunol. Methods 35:1-21 (1980); Kohler and Milstein, Nature 256:495-497 (1975)), the trioma technique, the human B-cell hybridoma technique (Kozbor et al., Immunology Today 4:72 (1983); Cole et al., in Monoclonal Antibodies and Cancer Therapy, Alan R. Liss, Inc. (1985), pp. 77-96). Antibodies may also be generated by the known techniques of phage display and in vitro immunization.

Polyclonal Antibodies

5

10

15

20

25

30

35

Polyclonal antibodies are heterogeneous populations of antibody molecules derived from the sera of animals immunized with an antigen, such as an inventive protein or an antigenic derivative thereof.

Polyclonal antiserum, containing antibodies to heterogeneous epitopes of a single protein, can be prepared by immunizing suitable animals with the expressed protein described above, which can be unmodified or modified, as known in the art, to enhance immunogenicity. Immunization methods include subcutaneous or intraperitoneal injection of the polypeptide.

Effective polyclonal antibody production is affected by many factors related both to the antigen and to the host species. For example, small molecules tend to be less immunogenic than other and may require the use of carriers and/or adjuvant. In addition, host animal response may vary with site of inoculation. Both inadequate or excessive doses of antigen may result in low titer antisera. In general, however, small doses (high ng to low µg levels) of antigen administered at multiple intradermal sites appears to be most reliable. Host animals may include but are not limited to rabbits, mice, chickens and rats, to name but a few. An effective immunization protocol for rabbits can be found in Vaitukaitis, J. et al., J. Clin. Endocrinol. Metab. 33:988-991 (1971).

The protein immunogen may be modified or administered in an adjuvant in order to increase the protein's antigenicity. Methods of increasing the antigenicity of a protein are well known in the art and include, but are not limited to coupling the antigen with a heterologous protein (such as globulin β -galactosidase) or through the inclusion of an adjuvant during immunization. Adjuvants include Freund's (complete and incomplete), mineral gels such as aluminum hydroxide, surface active substances such as lysolecithin pluronic polyols, polyanions, peptides oil emulsions, keyhole limpet hemocyanin dinitrophenol potentially useful human adjuvants such as BCG (bacille Calmette-Guerin) and Corynebacterium parvum.

Booster injections can be given at regular intervals, with at least one usually being required for optimal antibody production.

The antiserum may be harvested when the antibody titer begins to fall. Titer may be determined semi-quantitatively, for example, by double immunodiffusion in agar against known concentrations of the antigen. See, for example, Ouchterlony et al., Chap. 19 in: Handbook of Experimental Immunology, Wier, ed, Blackwell (1973). Plateau concentration of antibody is usually in the range of [].] to 0.2 mg/ml of serum (about 12 μ M). The antiserum may be purified by affinity chromatography using the immobilized immunogen carried on a solid support. Such methods of affinity chromatography are well known in the art.

Affinity of the antisera for the antigen may be determined by preparing competitive binding curves, as described, for example, by Fisher, Chap. 42 in: Manual of Clinical Immunology, second edition, Rose and Friedman, eds., Amer. Soc. For Microbiology, 15 Washington, D.C. (1980).

In addition to using protein an the immunogen, DNA molecules may be used directly. In this manner, a DNA encoding the protein immunogen is administered. Boosting and harvesting is done in a manner analogous to that detailed above. Yet another method of producing antibodies entails immunizing chickens and harvesting the antibodies from their eggs.

Monoclonal Antibodies

10

20

25

Monoclonal antibodies (MAbs), are homogeneous populations of antibodies to a particular antigen. They may be obtained by any technique which provides for the production of antibody molecules by continuous cell lines in culture or *in vivo*. MAbs may be produced by making hybridomas which are immortalized cells capable of secreting a specific monoclonal antibody.

Monoclonal antibodies to any of the proteins: peptides and epitopes thereof described herein can be prepared from murine hybridomas according to the classical method of Kohler: G. and Milstein: C.: Nature 256:495-497 (1975) (and U.S. Patent No. 5 4:376:110) or modifications of the methods thereof: such as the human B-cell hybridoma technique (Kosbor et al.: 1983: Immunology Today 4:72: Cole et al.: 1983: Proc. Natl. Acad. Sci. USA 80: 2026-2030): and the EBV-hybridoma technique (Cole et al.: 1985: MONOCLONAL ANTIBODIES AND CANCER THERAPY: Alan R. Liss: Inc.: pp. 10 77-96).

In one method a mouse is repetitively inoculated with a few micrograms of the selected protein over a period of a few weeks. The mouse is then sacrificed, and the antibody producing cells of the spleen are isolated.

15

20

25

30

35

The spleen cells are fused, typically using polyethylene glycol, with mouse myeloma cells, such as SP2/D-Agl4 myeloma cells. The excess, unfused cells are destroyed by growth of the system on selective media comprising aminopterin (HAT media). The successfully fused cells are diluted, and aliquots are plated to microliter plates where growth is continued. Antibody—producing clones (hybridomas) are identified by detection of antibody in the supernatant fluid of the wells by immunoassay procedures. These include ELISA, as originally described by Engvall, Meth. Enzymol. 70:419 (1980), western blot analysis, radioimmunoassay (Lutz et al., Exp. Cell Res. 175:109-124 (1988)) and modified methods thereof.

Selected positive clones can be expanded and their monoclonal antibody product harvested for use. Detailed procedures for monoclonal antibody production are described in Davis, L. et al. BASIC METHODS IN MOLECULAR BIOLOGY, Elsevier, New York. Section 21-2 (1989). The hybridoma clones may be cultivated in vitro or in vivo, for instance as ascites. Production of high titers of mAbs in vivo makes this the presently preferred method of production. Alternatively, hybridoma culture in hollow fiber bioreactors provides a continuous high yield source of monoclonal antibodies.

The antibody class and subclass may be determined using procedures known in the art (Campbell, A.M., Monoclonal Antibody

Technology: Laboratory Techniques in Biochemistry and Molecular Biology: Elsevier Science Publishers: Amsterdam: The Netherlands (1984)). MAbs may be of any immunoglobulin class including IgG: IgM: IgE: IgA: IgD and any subclass thereof. Methods of purifying monoclonal antibodies are well known in the art.

Antibody Derivatives and Fragments

10

15

20

25

30

35

Fragments or derivatives of antibodies include any portion of the antibody which is capable of binding the target antigen, or a specific portion thereof. Antibody derivatives include polyspecific (e.g., bi-specific) antibodies, which contain binding sites specific for two or more different epitopes. These epitopes may be from the same or different inventive molecules or one or more epitope may be from a molecule not specifically disclosed here.

Antibody fragments specifically include F(ab')₂₁ Fab₁ Fab' and Fv fragments. These can be generated from any class of antibody: but typically are made from IgG or IgM. They may be made by conventional recombinant DNA techniques or: using the classical method: by proteolytic digestion with papain or pepsin. See CURRENT PROTOCOLS IN IMMUNOLOGY: chapter 2: Coligan et al.: eds:: (John Wiley & Sons 1991-92).

F(ab')2 fragments are typically about 110 kDa (IgG) or about 150 kDa (IgM) and contain two antigen-binding regions; joined at the hinge by disulfide bond(s). Virtually all; if not all; of the Fc is absent in these fragments. Fab' fragments are typically about 55 kDa (IgG) or about 75 kDa (IgM) and can be formed; for example; by reducing the disulfide bond(s) of an $F(ab')_2$ fragment. The resulting free sulfhydryl group(s) may be used to conveniently conjugate Fab' fragments to other molecules; such as detection reagents (e.g., enzymes).

Fab fragments are monovalent and usually are about 50 kDa (from any source). Fab fragments include the light (L) and heavy (H) chain variable (V_L and V_{H} respectively) and constant (C_L C_{H} respectively) regions of the antigen-binding portion of the antibody. The H and L portions are linked by an intramolecular disulfide bridge.

Fv fragments are typically about 25 kDa (regardless of source) and contain the variable regions of both the light and

heavy chains (V_L and V_H , respectively). Usually, the V_L and V_H chains are held together only by non-covalent interacts and, thus, they readily dissociate. They do, however, have the advantage of small size and they retain the same binding properties of the larger Fab fragments. Accordingly, methods have been developed to crosslink the V_L and V_H chains, using, for example, glutaraldehyde (or other chemical crosslinkers), intermolecular disulfide bonds (by incorporation of cysteines) and peptide linkers. The resulting Fv is now a single chain (i.e., SCFv).

Other antibody derivatives include single chain antibodies (U.S. Patent 4.946.778; Bird, Science 242:423-426 (1988); Huston et al., Proc. Natl. Acad. Sci. USA 85:5879-5883 (1988); and Ward et al., Nature 334:544-546 (1989)). Single chain antibodies are formed by linking the heavy and light chain fragments of the Fv region via an aminomacid bridge, resulting in a single chain FV (SCFv).

10

15

20

25

35

One preferred method involves the generation of scFvs by recombinant methods, which allows the generation of Fvs with new specificities by mixing and matching variable chains different antibody sources. In a typical method, a recombinant would be provided which comprises the appropriate regulatory elements driving expression of a cassette region. cassette region would contain a DNA encoding a peptide linker, with convenient sites at both the 5' and 3' ends of the linker for generating fusion proteins. The DNA encoding a variable region(s) of interest may be cloned in the vector to form fusion proteins with the linker, thus generating an scfv.

In an exemplary alternative approach, DNAs encoding two Fvs may be ligated to the DNA encoding the linker, and the resulting 30 tripartite fusion may be ligated directly into a conventional expression vector. The scFv DNAs generated any of these methods may be expressed in prokaryotic or eukaryotic cells, depending on the vector chosen.

Antibody fragments which recognize specific epitopes may be generated by known techniques. For example, such fragments include but are not limited to: the $F(ab')_2$ fragments which can be produced by pepsin digestion of the antibody molecule and the Fab fragments which can be generated by reducing the disulfide bridges

-569-

of the F(ab)₂ fragments. Alternatively, Fab expression libraries may be constructed (Huse et al., 1989, *Science*, 246:1275-1281) to allow rapid and easy identification of monoclonal Fab fragments with the desired specificity.

Derivatives also include "chimeric antibodies" (Morrison et al., Proc. Natl. Acad. Sci., 81:6851-6855 (1984); Neuberger et al., Nature, 312:604-608 (1984); Takeda et al., Nature, 314:452-454 (1985)). These chimeras are made by splicing the DNA encoding a mouse antibody molecule of appropriate specificity with, for instance, DNA encoding a human antibody molecule of appropriate specificity. Thus, a chimeric antibody is a molecule in which different portions are derived from different animal species, such as those having a variable region derived from a murine mAb and a human immunoglobulin constant region. These are also known sometimes as "humanized" antibodies and they offer the added advantage of at least partial shielding from the human immune system. They are, therefore, particularly useful in therapeutic in vivo applications.

Labeled Antibodies

10

35

20 The present invention further provides the above-described antibodies in detectably labeled form. Antibodies can be detectably labelled through the use of radioisotopes, affinity labels (such as biotin, avidin, etc.), enzymatic labels (such as horseradish peroxidase, alkaline phosphatase, etc.) fluorescent 25 labels (such as FITC or rhodamine, etc.), paramagnetic atoms, etc. Procedures for accomplishing such labeling are well-known in the art, for example see (Sternberger et al., J. Histochem. Cytochem. 18:315 (1970); Bayer et al., Meth. Enzym. L2:308 (1979); Engval et al., Immunol. 109:129 (1972); Goding, J. Immunol. Meth. 13:215 30 The labeled antibodies of the present invention can be used for in vitro, in vivo, and in situ diagnostic assays.

Immobilized Antibodies

The foregoing antibodies also may be immobilized on a solid support. Examples of such solid supports include plastics such as polycarbonate, complex carbohydrates such as agarose and sepharose, acrylic resins and such as polyacrylamide and latex beads. Techniques for coupling antibodies to such solid supports

-570-

are well known in the art (Weir et al., "Handbook of Experimental Immunology" 4th Ed., Blackwell Scientific Publications, Oxford, England, Chapter 10 (1986); Jacoby et al., Meth. Enzym. 34 Academic Press, N.Y. (1974)). The immobilized antibodies of the present invention can be used for in vitro, in vivo, and in situ assays as well as for immunoaffinity purification of the proteins of the present invention.

THERAPEUTIC AND DIAGNOSTIC COMPOSITIONS

10

15

20

25

30

35

The proteins, antibodies and polynucleotides of the present invention can be formulated according to known methods to prepare pharmaceutically useful compositions, whereby these materials, or their functional derivatives, are combined in admixture with a pharmaceutically acceptable carrier vehicle. Suitable vehicles and their formulation, inclusive of other human proteins, e.g., human serum albumin, are described, for example, in Remington's Pharmaceutical Sciences (16th ed., Osol, A., Ed., Mack, Easton PA In order to form a pharmaceutically acceptable (1980)). for effective composition suitable administration. compositions will contain an effective amount of one or more of the agents of the present invention, together with a suitable amount of carrier vehicle.

Pharmaceutical compositions for use in accordance with the present invention may be formulated in conventional manner using one or more physiologically acceptable carriers or excipients. Thus, the compounds and their physiologically acceptable salts and solvate may be formulated for administration by inhalation or insufflation (either through the mouth or the nose) or oral, buccal, parenteral or rectal administration.

For oral administration, the pharmaceutical compositions may take the form of, for example, tablets or capsules prepared by conventional means with pharmaceutically acceptable excipients such as binding agents (e.g., pregelatinised maize starch, polyvinylpyrrolidone or hydroxypropyl methylcellulose); fillers <math>(e.g., lactose, microcrystalline cellulose or calcium hydrogen phosphate); lubricants <math>(e.g., magnesium stearate, talc or silica); disintegrants (e.g., potato starch or sodium starch glycolate); or

wetting agents (e.g., sodium lauryl sulphate). The tablets may be coated by methods well known in the art. Liquid preparations for oral administration may take the form of, for example, solutions, syrups or suspensions, or they maybe presented as a dry product for constitution with water or other suitable vehicle before use. Such liquid preparations may be prepared by conventional means with pharmaceutically acceptable additives such as suspending sorbitol (e.g., syrupı cellulose derivatives hydrogenated edible fats); emulsifying agents (e.g., lecithin or acacia); non-aqueous vehicles (e.g., almond oil, oily esters, ethyl alcohol or fractionated vegetable oils); and preservatives (e.g., methyl or propyl-p-hydroxybenzoates or sorbic acid). The preparations may also contain buffer salts, flavoring, coloring and sweetening agents as appropriate.

10

15

20

25

30

35

Preparations for oral administration may be suitably formulated to give controlled release of the active compound. For buccal administration the composition may take the form of tablets or lozenges formulated in conventional manner.

For administration by inhalation, the compounds for use according to the present invention are conveniently delivered in the form of an aerosol spray presentation from pressurized packs or a nebuliser, with the use of a suitable propellant, e.g., dichlorodifluoromethane, trichlorofluoromethane, dichlorotetrafluoroethane, carbon dioxide or other suitable gas. In the case of a pressurized aerosol the dosage unit may be determined by providing a valve to deliver a metered amount. Capsules and cartridges of, e.g. gelatin for use in an inhaler or insufflator may be formulated containing a powder mix of the compound and a suitable powder base such as lactose or starch.

The compounds may be formulated for parenteral administration by injection, e.g., by bolus injection or continuous infusion. Formulations for injection may be presented in unit dosage form, e.g., in ampules or in multi-dose containers, with an added preservative. The compositions may take such forms as suspensions, solutions or emulsions in oily or aqueous vehicles, and may contain formulatory agents such as suspending, stabilizing and/or dispersing agents. Alternatively, the active ingredient

may be in powder form for constitution with a suitable vehicle e.g., sterile pyrogen-free water, before use.

The compounds may also be formulated in rectal compositions such as suppositories or retention enemas, e.g., containing 5 conventional suppository bases such as cocoa butter or other glycerides.

In addition to the formulations described previously, the compounds may also be formulated as a depot preparation. Such long acting formulations may be administered by implantation (for example subcutaneously or intramuscularly) or by intramuscular injection. Thus, for example, the compounds may be formulated with suitable polymeric or hydrophobic materials (for example as an emulsion in an acceptable oil) or ion exchange resins, or as sparingly soluble derivatives, for example, as a sparingly soluble salt.

The compositions may, if desired, be presented in a pack or dispenser device which may contain one or more unit dosage forms containing the active ingredient. The pack may for example comprise metal or plastic foil, such as a blister pack. The pack or dispenser device may be accompanied by instructions for administration.

RECOMBINANT CONSTRUCTS AND EXPRESSION

10

15

20

25

35

The present invention further provides recombinant DNA constructs comprising one or more of the nucleotide sequences of the present invention. The recombinant constructs of the present · invention comprise a vector, such as a plasmid or viral vector, into which a DNA or DNA fragment, typically bearing an open reading frame, is inserted, in either orientation. The gene products encoded by the subject DNAs may be produced by recombinant DNA technology using techniques well known in the art. See, for example, the techniques described in Sambrook et al., 1989, supra, and Ausubel et al., 1989, supra. Alternatively, the DNA sequences may be chemically synthesized using, for example, See, for example, the techniques described in OLIGONUCLEOTIDE SYNTHESIS, 1984, Gait, ed., IRL Press, Oxford, which is incorporated by reference herein in its entirety. may be assembled from fragments and short oligonucleotide linkers,

or from a series of oligonucleotides. The are preferably made by RT-PCR methods. The resulting synthetic gene is capable of being expressed in a recombinant vector.

In some cases the recombinant constructs will be expression vectors, which are capable of expressing the RNA and/or protein products of the encoded DNA(s). Thus, the vector may further comprise regulatory sequences, including for example, a promoter, operably linked to the open reading frame (ORF). The vector may further comprise a selectable marker sequence.

10

15

30

35

Specific initiation signals may also be required for efficient translation of inserted target gene coding sequences. These signals include the ATG initiation codon and adjacent In cases where a target DNA includes sequences. its initiation codon and adjacent sequences is inserted into the appropriate expression vector, no additional translation control signals may be needed. However, in cases where only a portion of is used, exogenous translational control signals, including, perhaps, the ATG initiation codon, must be provided. Furthermore, the initiation codon must be in phase with the 20 reading frame of the desired coding sequence to ensure translation of the entire target. These exogenous translational control signals and initiation codons can be of a variety of origins, both natural and synthetic. The efficiency of expression may be enhanced by the inclusion of appropriate transcription enhancer elements, transcription terminators, etc. (see Bittner et al., 25 Methods in Enzymol. 153:516-544 (1987)). Some appropriate cloning and expression vectors for use with prokaryotic and eukaryotic hosts are described by Sambrook, et al., in Molecular Cloning: A Laboratory Manual, Second Edition, Cold Spring Harbor, New York (1989), the disclosure of which is hereby incorporated by reference.

If desired, to enhance expression and facilitate proper protein folding, the codon context and codon pairing of the sequence may be optimized for the particular expression organism, as explained by Hatfield et al., U.S. Patent No. 5,082,767.

The present invention further provides host cells containing at least one of the DNAs of the present invention. The host cell can be virtually any cell for which expression vectors are

available. It may be for example a higher eukaryotic host cells such as a mammalian cells a lower eukaryotic host cells such as a yeast cells or the host cell can be a prokaryotic cells such as a bacterial cell. Introduction of the recombinant construct into the host cell can be effected by calcium phosphate transfections DEAEs dextran mediated transfections or electroporation (Davis et al., Basic Methods in Molecular Biology (1986)).

A wide variety of expression systems are available, such as: yeast (e.g. Saccharomyces, Pichia) transformed with recombinant yeast expression vectors containing the target DNA; insect cell systems infected with recombinant virus expression vectors (e.g., baculovirus) containing the target DNA sequences; plant cell systems infected with recombinant virus expression vectors (e.g., cauliflower mosaic virus, CaMV; tobacco mosaic virus, TMV) or transformed with recombinant plasmid expression vectors (e.g. Ti plasmid) containing target DNA coding sequences; or mammalian cell systems (e.g. COS, CHO, BHK, 293, 3T3) harboring recombinant expression constructs containing promoters derived from the genome of mammalian cells (e.g., metallothionein promoter) or from mammalian viruses (e.g., the adenovirus late promoter; the vaccinia virus 7.5K promoter).

Depending on the system chosen, the resulting product may differ. For example, proteins expressed in most bacterial cultures, e.g., E. coli, will be free of glycosylation modifications; polypeptides or proteins expressed in yeast will have a glycosylation pattern different from that expressed in mammalian cells.

Vectors

10

15

20

25

30

35

Generally, recombinant expression vectors will include origins of replication and selectable markers permitting selection of the host cell, e.g., the ampicillin resistance gene of E. coli and S. cerevisiae TRPL gene, and a promoter derived from a highly-expressed gene to direct transcription of a downstream structural sequence. Such promoters can be derived from operons encoding glycolytic enzymes such as 3-phosphoglycerate kinase (PGK), α -factor, acid phosphatase, or heat shock proteins, among others. The heterologous structural sequence is assembled in appropriate

-575-



phase with translation initiation and termination sequence, and in one aspect of the invention, a leader sequence capable of directing secretion of translated protein into the periplasmic space or extracellular medium. Optionally, the heterologous sequence can encode a fusion protein including an N-terminal or C-terminal identification peptide imparting desired characteristics, e.g., stabilization or simplified purification of expressed recombinant product.

Bacterial Expression

10

15

20

25

30

35

Useful expression vectors for bacterial use are constructed by inserting a structural DNA sequence encoding a desired protein together with suitable translation initiation and termination signals in operable reading phase with a functional promoter. The vector will comprise one or more phenotypic selectable markers and an origin of replication to ensure maintenance of the vector and if desirable, to provide amplification within the host. Suitable prokaryotic hosts for transformation include E. coli, Bacillus subtilis, Salmonella typhimurium and various species within the genera Pseudomonas, Streptomyces, and Staphylococcus, although others may, also be employed as a matter of choice.

Bacterial vectors may be for example bacteriophage plasmid or cosmid-based. These vectors can comprise a selectable marker and bacterial origin of replication derived from commercially available plasmids typically containing elements of the well known cloning vector pBR322 (ATCC 37017). Such commercial vectors include for example GEM 1 (Promega Biotec Madison, WI, USA), pBs, phagescript, PsiX174, pBluescript SK, pBs KS, pNH&a, pNH1&a, pNH4&a (Stratagene); pTrc99A, pKK223-3, pKK233-3, pKK232-8, pDR540, and pRIT5 (Pharmacia).

These "backbone" sections are combined with an appropriate promoter and the structural sequence to be expressed. Bacterial promoters include lac, T3, T7, lambda P_R or P_L , trp, and ara.

Following transformation of a suitable host strain and growth of the host strain to an appropriate cell density, the selected promoter is derepressed/induced by appropriate means (e.g., temperature shift or chemical induction) and cells are cultured for an additional period. Cells are typically harvested by

centrifugation, disrupted by physical or chemical means, and the resulting crude extract retained for further purification.

In bacterial systems, a number of expression vectors may be advantageously selected depending upon the use intended for the protein being expressed. For example, when a large quantity of such a protein is to be produced, for the generation of antibodies or to screen peptide libraries, for example, vectors which direct the expression of high levels of fusion protein products that are readily purified may be desirable. Such vectors include, but are not limited, to the *E. coli* expression vector pUR278 (Ruther et al., 1983, *EMBO J.* 2:1791), in which the coding sequence may be ligated into the vector in frame with the lac Z coding region so that a fusion protein is produced; pIN vectors (Inouye et al., 1985, Nucleic Acids Res., 13:3101-3109; Van Heeke et al., 1989, J. Biol. Chem., 264:5503-5509); pET vectors, Studier et al., Methods in Enzymology 185: 60-89 (Academic Press 1990); and the like.

Moreover, pGEX vectors may be used to express foreign polypeptides as fusion proteins with glutathione S-transferase (GST). In general, such fusion proteins are soluble and easily can be purified from lysed cells by adsorption to glutathione-agarose beads followed by elution in the presence of free glutathione. The pGEX vectors are designed to include thrombin or factor Xa protease cleavage sites so that the cloned target gene protein can be released from the GST moiety.

In a one embodiment, full length cDNA sequences are appended with in-frame BamHI sites at the amino terminus and EcoRI sites at the carboxyl terminus using standard PCR methodologies (Innis et al., 1990, supra) and ligated into the pGEX-2TK vector (Pharmacia, Uppsala, Sweden). The resulting cDNA construct contains a kinase recognition site at the amino terminus for radioactive labeling and glutathione S-transferase sequences at the carboxyl terminus for affinity purification (Nilsson, et al. 1985, EMBO J. 4: 1075; Zabeau and Stanley, 1982, EMBO J. 1:1217.

Eukaryotic Expression

20

25

Various mammalian cell culture systems can also be employed to express recombinant protein. Examples of mammalian expression systems include the COS-7 lines of monkey kidney fibroblasts.

described by Gluzman, Cell 23:175 (1981), and other cell lines capable of expressing a compatible vector, for example, the Cl27, 3T3, CHO, HeLa and BHK cell lines. Mammalian expression vectors will comprise an origin of replication, a suitable promoter and also any necessary ribosome binding polyadenylation site, splice donor and acceptor sites. transcriptional termination sequences and 51 flanking nontranscribed sequences. DNA sequences derived from the SV40 viral genome, for example, SV4O origin, early promoter, enhancer, splice, and polyadenylation sites may be used to provide the required nontranscribed genetic elements.

10

15

20

25

30

Mammalian promoters include CMV immediate early, HSV thymidine kinase, early and late SV40, LTRs from retrovirus, and mouse metallothionein-I. Exemplary mammalian vectors include pWLneo, pSV2cat, p0G44, pXT1, pSG (Stratagene) pSVK3, pBPV, pMSG, and pSVL (Pharmacia). Selectable markers include CAT (chloramphenicol transferase).

In mammalian host cells, a number of viral-based expression systems may be utilized. In cases where an adenovirus is used as an expression vector, the coding sequence of interest may be ligated to an adenovirus transcription/translation control complex, e.g., the late promoter and tripartite leader sequence. This chimeric gene may then be inserted in the adenovirus genome by in vitro or in vivo recombination. Insertion in a non-essential region of the viral genome (e.g., region El or E3) will result in a recombinant virus that is viable and capable of expressing a target protein in infected hosts. (E.g., See Logan et al., 1984, Proc. Natl. Acad. Sci. USA 81:3655-3659).

In one embodiment, cDNA sequences encoding the full-length open reading frames are ligated into pCMVR replacing the R-galactosidase gene such that cDNA expression is driven by the CMV promoter (Alam, 1990, Anal. Biochem. 188: 245-254; MacGregor et al., 1989, Nucl. Acids Res. 17: 2365; Norton et al. 1985, Mol. Cell. Biol. 5: 281).

In addition, a host cell strain may be chosen which modulates the expression of the inserted sequences, or modifies and processes the gene product in the specific fashion desired. Such modifications (e.g., glycosylation) and processing (e.g.,

5

10

15

20

25

30

35

cleavage) of protein products may be important for the function of the protein. Different host cells have characteristic and specific mechanisms for the post-translational processing and modification of proteins.

Appropriate cell lines or host systems can be chosen to ensure the correct modification and processing of the foreign protein expressed. To this end, eukaryotic host cells which possess the cellular machinery for proper processing of the primary transcript, glycosylation, and phosphorylation of the gene product may be used. Such mammalian host cells include but are not limited to CHO, VERO, BHK, HeLa, COS, MDCK, 293, 3T3, WI38, etc.

For long-term, high-yield production of recombinant proteins in eukaryotic cells, stable expression is preferred. Rather than using expression vectors which contain viral origins of replication, host cells can be transformed with DNA controlled by appropriate expression control elements (e.g., promoter, enhancer, sequences, transcription terminators, polyadenylation sites, etc.), and a selectable marker.

Following the introduction of the foreign DNA, engineered cells may be allowed to grow for 1-2 days in an enriched media, and then are switched to a selective media. The selectable marker in the recombinant plasmid confers resistance to the selection and allows cells to stably integrate the plasmid into their chromosomes and grow to form foci which in turn can be cloned and expanded into cell lines. This method may advantageously be used to engineer cell lines which express the target protein. Such engineered cell lines may be particularly useful in screening and evaluation of compounds that affect the endogenous activity of the protein.

A number of selection systems may be used, including but not limited to the herpes simplex virus thymidine kinase (Wigler, et al., Cell 11:223 (1977)), hypoxanthine-guanine phosphoribosyltransferase (Szybalska et al., Proc. Natl. Acad. Sci. USA 48:2026 (1962)), and adenine phosphoribosyltransferase (Lowy, et al., Cell 22:817 (1980)) genes can be employed in the hypothesis of april cells, respectively. Also, antimetabolite resistance can be used as the basis of selection for dhfr, which

confers resistance to methotrexate (Wigler, et al., Proc. Natl. Acad, Sci. USA 77:35b7 (1980)); O'Hare, et al., 1981, Proc. Natl. Acad. Sci. USA 78:1527); gpt, which confers resistance to mycophenolic acid (Mulligan et al., Proc. Natl. Acad. Sci. USA 78:2072 (1981)); neo, which confers resistance to the aminoglycoside G-418 (Colberre-Garapin, et al., 1981, J. Mol. Biol. 150:1); and hydro, which confers resistance to hygromycin (Santerre, et al., 1984, Gene 30:147) genes.

An alternative fusion protein system allows for the ready purification of non-denatured fusion proteins expressed in human cell lines (Janknecht, et al., Proc. Natl. Acad. Sci. USA &&: 8972-8976 (1991)). In this system, the gene of interest is subcloned into a vaccinia-based plasmid such that the gene's open reading frame is translationally fused to an amino-terminal tag consisting of six histidine residues. Extracts from cells infected with recombinant vaccinia virus are loaded onto Ni²+ nitriloacetic acid-agarose columns and histidine-tagged proteins are selectively eluted with imidazole-containing buffers.

10

15

20

25

30

35

insect systema Autographa californica polyhedrosis virus (AcNPV) is used as a vector to express foreign The virus grows in Spodoptera frugiperda cells. target coding sequence may be cloned individually into nonessential regions (for example the polyhedrin gene) of the virus and placed under control of an AcNPV promoter (for example the polyhedrin promoter). Successful insertion of a target gene coding sequence will result in inactivation of the polyhedrin gene and production of non-occluded recombinant virus (i.e., virus lacking the proteinaceous coat coded for by the polyhedrin gene). These recombinant viruses are then used to infect Spodoptera frugiperda cells in which the inserted gene is expressed. (E.g., see Smith et al., 1983, J. Virol. 46: 584; Smith, U.S. Patent No. 4,235,051).

While the present proteins can be expressed in recombinant systems, as described above, cell-free translation systems can also be employed to produce such proteins using RNAs derived from the DNA constructs of the present invention.

Purification of Recombinant Proteins

10

15

20

25

30

35

Recombinant proteins produced may be isolated by host cell lysis. This may be followed by one or more salting-out, aqueous ion exchange or size exclusion chromatography steps. Finally, high performance liquid chromatography (HPLC) can be employed for final purification steps. Microbial cells employed in expression of proteins can be disrupted by any convenient method, including freeze-thaw cycling, sonication, mechanical disruption, or use of cell lysing agents, like lysozyme and chelators.

If inclusion bodies are formed in bacterial systems, they may be extracted from cell pellets using, for example, detergents, reducing agents, salts, urea, guanidinium chloride and extremes of pH (e.g. < 4 or > 10). If denaturation occurs, protein refolding steps (e.g., dialysis) can be used, as necessary, in completing configuration of the mature protein. If disulfide bridges are present in the native protein, they may be reoxidized using known methods.

By way of specific non-limiting example, the recombinant bacterial cells, for example E. coli, are grown in any of a number of suitable media, for example LB, and the expression of the recombinant protein induced by adding IPTG (e.g., lac operatorpromoter) to the media or switching incubation to a higher temperature (e.g., λ cI⁸⁵⁷). After culturing the bacteria for a further period of between 2 and 24 hours, the cells are collected by centrifugation and washed to remove residual mediabacterial cells are then lysed, for example, by disruption in a cell homogenizer and centrifuged to separate the cell membranes from the soluble cell components. If the protein aggregates into inclusion bodies, this centrifugation can be performed under conditions whereby the dense inclusion bodies are selectively enriched by incorporation of sugars such as sucrose into the buffer and centrifugation at a selective speed. bodies can then be washed in any of several solutions to remove some of the contaminating host proteins, then solubilized in solutions containing high concentrations of urea (e.g. chaotropic agents such as guanidinium hydrochloride presence of reducing agents such as ß-mercaptoethanol or DTT (dithiothreitol).

At this stage it may be advantageous to incubate the protein for several hours under conditions suitable for the protein to undergo a refolding process into a conformation which more closely resembles that of the native protein. Such conditions generally 5 include low protein concentrations less than 500 μg/ml), levels of reducing agent, concentrations of urea less than 2 M and often the presence of reagents such as a mixture of reduced and oxidized glutathione which facilitate the interchange disulphide bonds within the protein molecule. The refolding process can be monitored, for example, by SDS-PAGE or with antibodies which are specific for the native molecule. Following refolding, the protein can then be purified further and separated from the refolding mixture by chromatography on any of several supports including ion exchange resins, gel permeation resins or on a variety of affinity columns.

Labeling Proteins

When used as a component in assay systems such as those described, below, the target protein may be labeled, either directly or indirectly, to facilitate detection of the present 20 res-like molecules either in vitro or in vivo. Any of a variety of suitable labeling systems may be used including but not limited to radioisotopes such as 1251; enzyme labeling systems that generate a detectable colorimetric signal or light when exposed to substrate; and fluorescent labels.

25 Where recombinant DNA technology is used for protein production the it may be advantageous to engineer fusion proteins that can facilitate labeling immobilization and/or detection. These fusion proteins may for example add amino acids which facilitate further chemical modification. They also may add a functional moiety such as an enzyme which directly facilitates detection.

-582-

PCT/IB01/02050

WO 01/98454 TRANSGENIC ANIMALS

10

20

25

30

35

The invention further contemplates animal models for studying the function of the present molecules and for overproducing the protein products. The disclosed DNA sequences may be used in conjunction with techniques for producing transgenic animals that are well known to those of skill in the art.

To prepare transgenic animals, target gene sequences may for example be introduced into, and overexpressed in, the genome of the animal of interest, or, if endogenous target gene sequences are present, they may either be overexpressed or, alternatively, be disrupted in order to underexpress or inactivate target gene expression, such as described for the disruption of apoE in mice (Plum et al., Cell 71: 343-353 (1992)).

In order to overexpress a target gene sequence, the coding portion of the target gene sequence may be ligated to a regulatory sequence which is capable of driving gene expression in the animal and cell type of interest. Such regulatory regions will be well known to those of skill in the art, and may be utilized in the absence of undue experimentation.

For underexpression of an endogenous target gene sequences such a sequence may be isolated and engineered such that when reintroduced into the genome of the animal of interests the endogenous target gene alleles will be inactivated. Preferably, the engineered target gene sequence is introduced via gene targeting such that the endogenous target sequence is disrupted upon integration of the engineered target gene sequence into the animal's genome. Animals of any species, including but not limited to mice rats rabbits guinea pigs pigs micro-pigs goats and non-human primates e.g., baboons monkeys and chimpanzees may be used to generate cardiovascular disease animal models. Goats cows and sheep are particularly preferred for producing protein in vivo.

Any technique known in the art may be used to introduce a target gene transgene into animals to produce the founder lines of transgenic animals. Such techniques include, but are not limited to pronuclear microinjection (Hoppe et al., U.S. Pat. No. 4-873-191 (1989)); retrovirus mediated gene transfer into germ lines (Van der Putten et al., Proc. Natl. Acad. Sci., USA 82:6148-

6152 (1985)); gene targeting in embryonic stem cells (Thompson et al., Cell 56:313-321 (1989)); electroporation of embryos (Lo, Mol. Cell. Biol. 3:1803-1814 (1983)); and sperm-mediated gene transfer (Lavitrano et al., Cell 57:717-723 (1989)); etc. For a review of 5 such techniques, see Gordon, Transgenic Animals, Intl. Rev. Cytol. 115:171-229 (1989).

The present invention provides for transgenic animals that carry the transgene in all their cells, as well as animals which carry the transgene in some, but not all their cells, i.e., mosaic animals. The transgene may be integrated as a single transgene or concatamers, e.g., head-to-head tandems or head-to-tail The transgene may also be selectively introduced into and activated in a particular cell type by following, for example, the teaching of Lasko et al., Proc. Natl. Acad. Sci. USA 89:3232-6236 (1992)). The regulatory sequences required for such a cell-type specific activation will depend upon the particular cell type of interest, and will be apparent to those of skill in the art. When it is desired that the target gene be integrated into the chromosomal site of the endogenous target . gene, gene targeting is preferred. Briefly, when such a technique is to be utilized, vectors containing some nucleotide sequences homologous to the endogenous target gene of interest are designed: for the purpose of integrating, via homologous recombination with chromosomal sequences, into and disrupting the function of the 25 nucleotide sequence of the endogenous target gene-

10

20

35

The transgene may also be selectively introduced into a particular cell type, thus inactivating the endogenous gene of interest in only that cell type, by following, for example, the teaching of Gu et al. Science 265: 103-106 (1994)). regulatory sequences required for such a cell-type specific inactivation will depend upon the particular cell interest, and will be apparent to those of skill in the art.

Once transgenic animals have been generated, the expression of the recombinant target gene and protein may be assayed utilizing standard techniques. Initial screening accomplished by Southern blot analysis or PCR techniques to analyze animal tissues to assay whether integration of transgene has taken place. The level of mRNA expression of the

PCT/IB01/02050 WO 01/98454

transgene in the tissues of the transgenic animals may also be assessed using techniques which include but are not limited to Northern blot analysis of tissue samples obtained from the animal, in situ hybridization analysis, and RT-PCR. Samples of target gene-expressing tissue, may also be evaluated immunocytochemically using antibodies specific for the target gene transgene gene product of interest.

The transgenic animals that express target gene mRNA or target gene transgene peptide (detected immunocytochemically, using antibodies directed against the target gene product's epitopes) at easily detectable levels should then be further evaluated to identify those animals which display characteristic susceptibility to increased carcinogenesis. Additionally specific cell types within the transgenic animals may be analyzed and assayed in vitro for cellular phenotypes characteristic of mutant phenotype.

10

15

20

25

35

Once target gene transgenic founder animals are produced. they may be bred, inbred, outbred, or crossbred to produce colonies of the particular animal. Examples of such breeding strategies include but are not limited to: outbreeding of founder animals with more than one integration site in order to establish separate lines; inbreeding of separate lines in order to produce compound target gene transgenics that express the target gene transgene of interest at higher levels because of the effects of additive expression of each target gene transgene; crossing of heterozygous transgenic animals to produce animals homozygous for a given integration site in order both to augment expression and eliminate the possible need for screening of animals by DNA analysis; crossing of separate homozygous lines to 30 compound heterozygous or homozygous lines; breeding animals to different inbred genetic backgrounds so as to examine effects of modifying alleles on expression of the target gene transgene and the possible development of carcinogenesis. One such approach is to cross the target gene transgenic founder animals with a wild type strain to produce an FL generation that exhibits increased susceptibility to carcinogenesis. The FL generation may then be inbred in order to develop a homozygous line, if it is found that homozygous target gene transgenic animals are viable.

Methods of generating "knockout" mice using homologous recombination in embryonic stem cells are well known in the art. Suitable methods are described for example in Mansour et al., Nature 335:348 (1988); Zijlstra et al., Nature 342:435 (1989) and 344:742 (1990); and Hasty et al., Nature 350:243 (1991). This genomic DNA can be obtained by conventional methods using the cDNA sequence as a probe in a commercially-available genomic DNA library.

Briefly, a genomic fragment is cleaved with a restriction endonuclease and a heterologous cassette containing a neomycin-resistance gene is inserted at the cleavage site. A suitable cassette is the GTI-II neo cassette described by Lufkin et al., Cell bb:lb05 (1991). The modified genomic fragment is cloned into a suitable targeting vector that is introduced into murine embryonic stem cells by electroporation. Cells that have undergone homologous recombination (and hence disruption of the gene) are selected by resistance to G418, and used to generate chimeric mice using well known methods. See Lufkin et al., supra. Traditional breeding methods then can be used to generate mice that are homozygous for the disrupted gene.

The phenotype of mice that are homozygous for the mutation then can be studied to provide insights into the role of the protein in, for example, carcinogenesis. These mice also can be used as models for developing new treatments for cancers. If this mutation is lethal in homozygous mice (for example during embryogenesis) heterozygous mice, which express only half the amount of the protein can also be studied.

GENE THERAPY APPLICATIONS

10

15

20

25

30

35

When mutations in the inventive protein, or in the elements controlling expression of that protein, are found to be associated with a malignant phenotype, control of cellular proliferation can be restored by gene therapy methods. For example, overexpression of the protein can be counteracted by concurrent expression of an antisense molecule that binds to and inhibits expression of the mRNA encoding the protein. Alternatively, overexpression can be inhibited in an analogous manner using a ribozyme that cleaves the mRNA. In another embodiment, where expression of a mutated

protein induces the malignant phenotype, concomitant expression of the non-mutated molecule via introduction of an exogenous gene may be used. Methods of using antisense and ribozyme technology to control gene expression, or of gene therapy methods for expression of an exogenous gene in this manner are well known in the art.

Each of these methods requires a system for introducing a vector into the cells containing the mutated gene. The vector encodes either an antisense or ribozyme transcript of the inventive protein. The construction of a suitable vector can be achieved by any of the methods well-known in the art for the insertion of exogenous DNA into a vector. See, e.g., Sambrook et al., Molecular Cloning (Cold Spring Harbor Press 2d ed. 1989), which is incorporated herein by reference. In addition, the prior art teaches various methods of introducing exogenous genes into cells in vivo. See Rosenberg et al., Science 242:1575-1578 (1988) and Wolff et al., PNAS 86:9011-9014 (1989), which are incorporated herein by reference. The routes of delivery include systemic administration and administration in situ-Well-known techniques include systemic administration with cationic liposomes, 20 administration in situ with viral vectors. Any one of the gene delivery methodologies described in the prior art is suitable for the introduction of a recombinant vector containing an inventive gene according to the invention into a MTX-resistant, transportdeficient cancer cell. A listing of present-day vectors suitable for the purpose of this invention is set forth in Hodgson, Bio/Technology 13: 222 (1995); which is incorporated by reference.

10

15

25

30

35

For example, liposome-mediated gene transfer is a suitable method for the introduction of a recombinant vector containing an inventive gene according to the invention into a MTX-resistant, transport-deficient cancer cell. The use of a cationic liposome, such as DC-Chol/DOPE liposome, has been widely documented as an appropriate vehicle to deliver DNA to a wide range of tissues through intravenous injection of DNA/cationic liposome complexes. See Caplen et al., Nature Med. 1:39-46 (1995) and Zhu et al., Science 261:209-211 (1993), which are herein incorporated by reference. Liposomes transfer genes to the target cells by fusing with the plasma membrane. The entry process is relatively efficient, but once inside the cell, the liposome-DNA complex has

no inherent mechanism to deliver the DNA to the nucleus. As such the most of the lipid and DNA gets shunted to cytoplasmic waste systems and destroyed. The obvious advantage of liposomes as a gene therapy vector is that liposomes contain no proteins, which thus minimizes the potential of host immune responses.

As another example, viral vector-mediated gene transfer is also a suitable method for the introduction of the vector into a target cell. Appropriate viral vectors include adenovirus vectors and adeno-associated virus vectors, retrovirus vectors and herpesvirus vectors.

10

15

20

35

Adenoviruses are linear, double stranded DNA viruses complexed with core proteins and surrounded by capsid proteins. The common serotypes 2 and 5, which are not associated with any human malignancies, are typically the base vectors. By deleting parts of the virus genome and inserting the desired gene under the control of a constitutive viral promoter, the virus becomes a replication deficient vector capable of transferring the exogenous DNA to differentiated, non-proliferating cells. To enter cells, the adenovirus fibre interacts with specific receptors on the cell surface, and the adenovirus surface proteins interact with the surface integrins. The virus penton-cell interaction provides the signal that brings the exogenous genecontaining virus into a cytoplasmic endosome. The adenovirus breaks out of the endosome and moves to the nucleus, the viral capsid falls apart, and the exogenous DNA enters the cell nucleus where it functions, in an epichromosomal fashion, to express the exogenous gene. Detailed discussions of the use of adenoviral vectors for gene therapy can be found in Berkner, Biotechniques 6:616-629 (1988) and Trapnell, Advanced Drug Delivery Rev. 12:185-(EPPE) incorporated which are herein by reference. Adenovirus-derived vectors. particularly non-replicative adenovirus vectors, are characterized by their ability to accommodate exogenous DNA of 7.5 kB, relative stability, wide host range, low pathogenicity in man, and high titers (104 to 105 plaque forming units per cell). See Stratford-Perricaudet et al., PNAS 89:2581 (1992).

Adeno-associated virus (AAV) vectors also can be used for the present invention. AAV is a linear single-stranded DNA parvovirus

that is endogenous to many mammalian species. AAV has a broad host range despite the limitation that AAV is a defective parvovirus which is dependent totally on either adenovirus or herpesvirus for its reproduction in vivo. The use of AAV as a vector for the introduction into target cells of exogenous DNA is well-known in the art. See, e.g., Lebkowski et al., Mole. & Cell. Biol. 8:3988 (1988), which is incorporated herein by reference. In these vectors, the capsid gene of AAV is replaced by a desired DNA fragment, and transcomplementation of the deleted capsid function is used to create a recombinant virus stock. infection the recombinant virus uncoats in the nucleus integrates into the host genome.

10

15

20

25

35

Another suitable virus-based gene delivery mechanism retroviral vector-mediated gene transfer. In general, retroviral vectors are well-known in the art. See Breakfield et al., Mole. Neuro. Biol. 1:339 (1987) and Shih et al., in Vaccines 85: 177 (Cold Spring Harbor Press 1985). A variety of retroviral vectors and retroviral vector-producing cell lines can be used for the present invention. Appropriate retroviral vectors include Moloney Murine Leukemia Virus, spleen necrosis virus, and vectors derived from retroviruses such as Rous Sarcoma Virus, Harvey Sarcoma avian leukosis virus. human immunodeficiency virus, myeloproliferative sarcoma virus, and mammary tumor virus. vectors include replication-competent and replication-defective In addition, amphotropic and xenotropic retroviral vectors. retroviral vectors can be used. In carrying out the invention, retroviral vectors can be introduced to a tumor directly or in the form of free retroviral vector producing-cell lines. producer cells include fibroblasts, neurons, glial 30 keratinocytes, hepatocytes, connective tissue cells, ependymal cells, chromaffin cells. See Wolff et al., PNAS 84:3344 (1989).

Retroviral vectors generally are constructed such that the majority of its structural genes are deleted or replaced exogenous DNA of interest, and such that the likelihood is reduced that viral proteins will be expressed. See Bender et al., J. Virol. 61:1639 (1987) and Armento et al., J. Virol. 61:1647 (1987), which are herein incorporated by reference. To facilitate expression of the antisense or ribozyme molecule, of the inventive

protein, a retroviral vector employed in the present invention must integrate into the genome of the host cell genome, an event which occurs only in mitotically active cells. The necessity for host cell replication effectively limits retroviral gene expression to tumor cells, which are highly replicative, and to a few normal tissues. The normal tissue cells theoretically most likely to be transduced by a retroviral vector, therefore, are the endothelial cells that line the blood vessels that supply blood to the tumor. In addition, it is also possible that a retroviral vector would integrate into white blood cells both in the tumor or in the blood circulating through the tumor.

10

15

20

25

30

35

The spread of retroviral vector to normal tissues, however, is limited. The local administration to a tumor of a retroviral vector or retroviral vector producing cells will restrict vector propagation to the local region of the tumor, minimizing transduction, integration, expression and subsequent cytotoxic effect on surrounding cells that are mitotically active.

Both replicatively deficient and replicatively competent retroviral vectors can be used in the invention, subject to their respective advantages and disadvantages. For instance, for tumors that have spread regionally, such as lung cancers, the direct injection of cell lines that produce replication-deficient vectors may not deliver the vector to a large enough area to completely eradicate the tumor, since the vector will be released only form the original producer cells and their progeny, and diffusion is limited. Similar constraints apply to the application replication deficient vectors to tumors that grow slowly, such as human breast cancers which typically have doubling times of 3D days versus the 24 hours common among human gliomas. The much shortened survival-time of the producer cells, probably no more than 7-14 days in the absence of immunosuppression, limits to only a portion of their replicative cycle the exposure of the tumor cells to the retroviral vector.

The use of replication-defective retroviruses for treating tumors requires producer cells and is limited because each replication-defective retrovirus particle can enter only a single cell and cannot productively infect others thereafter. Because these replication-defective retroviruses cannot spread to other tumor cells, they would be unable to completely penetrate a deep.

PCT/IB01/02050 WO 01/98454

multilayered tumor in vivo. See Markert et al., Neurosurg. 77: The injection of replication-competent retroviral vector particles or a cell line that produces a replicationcompetent retroviral vector virus may prove to be a more effective therapeutic because a replication competent retroviral vector will establish a productive infection that will transduce cells as long as it persists. Moreover, replicatively competent retroviral vectors may follow the tumor as it metastasizes, carried along and propagated by transduced tumor cells. The risks for complications are greater, with replicatively competent vectors, however. vectors may pose a greater risk then replicatively deficient vectors of transducing normal tissues, for instance. The risks of undesired vector propagation for each type of cancer and affected body area can be weighed against the advantages in the situation replicatively competent verses replicatively deficient retroviral vector to determine an optimum treatment.

10

15

20

25

35

Both amphotropic and xenotropic retroviral vectors may be used in the invention. Amphotropic viruses have a very broad host range that includes most or all mammalian cells, as is well known Xenotropic viruses can infect all mammalian cell to the art. cells except mouse cells. Thus, amphotropic and xenotropic retroviruses from many species, including cows, sheep, pigs, dogs, cats, rats, and mice, inter alia can be used to provide retroviral vectors in accordance with the invention, provided the vectors can transfer genes into proliferating human cells in vivo.

Clinical trials employing retroviral vector therapy treatment of cancer have been approved in the United States. See Culver, Clin. Chem. 40: 510 (1994). Retroviral vector-containing cells have been implanted into brain tumors growing in human patients. 30 See Oldfield et al., Hum. Gene Ther. 4: 39 (1993). retroviral vectors carried the HSV-1 thymidine kinase (HSV-tk) gene into the surrounding brain tumor cells, which conferred sensitivity of the tumor cells to the antiviral drug ganciclovir. Some of the limitations of current retroviral based cancer therapy, as described by Oldfield are: (1) the low titer of virus produced, (2) virus spread is limited to the region surrounding the producer cell implant, (3) possible immune response to the producer cell line, (4) possible insertional mutagenesis and

transformation of retroviral infected cells: (5) only a single treatment regimen of pro-drug: ganciclovir: is possible because the "suicide" product kills retrovirally infected cells and producer cells and (b) the bystander effect is limited to cells in direct contact with retrovirally transformed cells. See Bi et al.: Human Gene Therapy 4: 725 (1993).

Yet another suitable virus-based gene delivery mechanism is herpesvirus vector-mediated gene transfer. While much less is known about the use of herpesvirus vectors, replication-competent HSV-L viral vectors have been described in the context of antitumor therapy. See Martuza et al., Science 252: 854 (1991), which is incorporated herein by reference.

DIAGNOSTIC METHODS

10

15

20

25

30

35

The present invention also contemplates, for certain molecules described below, methods for diagnosis of human disease. In particular, patients can be screened for the occurrence of cancers, or likelihood of occurrence of cancers, associated with mutations in the encoded protein. DNA from tumor tissue obtained from patients suffering from cancer can be isolated and the gene encoding the protein can be sequenced. By examining a number of patients in this manner, mutations in the gene that are associated with a malignant cellular phenotype can be identified. addition, correlation of the nature of the observed mutations with subsequent observed clinical outcomes allows development prognostic model for the predicted outcome in a particular patient.

Screening for mutations conveniently can be carried out at the DNA level by use of PCR, although the skilled artisan will be aware that many other well known methods are available for the screening. PCR primers can be selected that flank known mutation sites, and the PCR products can be sequenced to detect the occurrence of the mutation. Alternatively, the 3' residue of one PCR primer can be selected to be a match only for the residue found in the unmutated gene. If the gene is mutated, there will be a mismatch at the 3' end of the primer, and primer extension cannot occur, and no PCR product will be obtained. Alternatively, primer mixtures can be used where the 3' residue of one primer is

-592-

any nucleotide other than the nonmutated residue. Observation of a PCR product then indicates that a mutation has occurred. Other methods of using, for example, oligonucleotide probes to screen for mutations are described, or example, in U.S. Patent No. 4,871,838, which is herein incorporated by reference in its entirety.

Alternatively, antibodies can be generated that selectively bind either mutated or non-mutated protein. The antibodies then can be used to screen tissue samples for occurrence of mutations in a manner analogous to the DNA-based methods described supra-

The diagnostic methods described above can be used not only for diagnosis and for prognosis of existing disease, but may also be used to predict the likelihood of the future occurrence of disease. For example, clinically healthy patients can be screened for mutations in the inventive molecule that correlate with later disease onset. Such mutations may be observed in the heterozygous state in healthy individuals. In such cases a single mutation event can effectively disable proper functioning of the gene and induce a transformed or malignant phenotype. This screening also may be carried out prenatally or neonatally.

DNA molecules according to the invention also are well suited for use in so-called "gene chip" diagnostic applications. Such applications have been developed by, inter alia, Synteni and Briefly, all or part of the DNA molecules of the Affymetrix. invention can be used either as a probe to screen a polynucleotide array on a "gene chip," or they may be immobilized on the chip used to identify other polynucleotides hybridization to the surface of the chip. In this manner, for example, related genes can be identified, or expression patterns of the gene in various tissues can be simultaneously studied. Such gene chips have particular application for diagnosis of disease, or in forensic analysis to detect the presence or absence of an analyte. Suitable chip technology is described for example, in Wodicka et al., Nature Biotechnology, 15:1359 (1997) which is hereby incorporated by reference in its entirety, and references cited therein.

-593-

PROTEIN-PROTEIN INTERACTIONS

10

20

25

30

35

Due to their similarity to certain known proteins, it is anticipated that some of the inventive protein molecules will interact with another class of cellular proteins. This is particularly true of those molecule containing leucine zipper motifs.

method suitable for Any detecting protein-protein interactions can be employed for identifying interacting targets. Among the traditional methods which can be employed are coimmunoprecipitation crosslinking and co-purification gradients or chromatographic columns. Utilizing procedures such as these allows for the identification of GAP gene products. Once identified, a GAP protein can be used, in conjunction with standard techniques, to identify its corresponding pathway gene. For example, at least a portion of the amino acid sequence of the pathway gene product can be ascertained using techniques well known to those of skill in the art, such as via the Edman degradation technique (see, e.g., Creighton, 1983, STRUCTURES AND MOLECULAR PRINCIPLES, W.H. Freeman & Co., N.Y., The amino acid sequence obtained can be used as a pp.34-49). guide for the generation of oligonucleotide mixtures that can be used to screen for pathway gene sequences. Screening can be accomplished, for example, by standard hybridization or PCR Techniques for the generation of oligonucleotide techniques. mixtures and for screening are well-known. (See e.g., Ausubel, supra, and PCR PROTOCOLS: A GUIDE TO METHODS AND APPLICATIONS, 1990: Innis et al., eds. Academic Press, Inc., New York).

10

15

20

25

30

Additionally, methods can be employed which result in the simultaneous identification of interacting target genes. One method which detects protein interactions in vivo, the two-hybrid system, is described in detail for illustration purposes only and not by way of limitation. One version of this system has been described (Chien et al., Proc. Natl. Acad. Sci. USA, 88: 9578-9582 (1991)) and is commercially available from Clontech (Palo Alto, CA).

Briefly, utilizing such a system, plasmids are constructed that encode two hybrid proteins: one consists of the DNA-binding domain of a transcription activator protein fused to a known protein, in this case an inventive protein, and the other contains

the activator protein's activation domain fused to an unknown protein (a putative GAP, for instance) that is encoded by a cDNA which has been recombined into this plasmid as part of a cDNA library. The plasmids are transformed into a strain of the yeast Saccharomyces cerevisiae that contains a reporter gene (e.g., lacZ) whose regulatory region contains the transcription activator's binding sites. Either hybrid protein alone cannot activate transcription of the reporter gene, the DNA-binding domain hybrid cannot because it does not provide activation function, and the activation domain hybrid cannot because it cannot localize to the activator's binding sites. Interaction of the two hybrid proteins reconstitutes the functional activator protein and results in expression of the reporter gene, which is detected by an assay for the reporter gene product.

15

20

25

30

35

The two-hybrid system or related methodology can be used to screen activation domain libraries for proteins that interact with a known "bait" gene product. By way of example, and not by way of limitation, gene products known to be involved in TH cell disorders subpopulation-related and/or differentiation. maintenance, and/or effector function of the subpopulations can be used as the bait gene products. Total genomic or cDNA sequences are fused to the DNA encoding on activation domain. This library and a plasmid encoding a hybrid of the bait gene product fused to the DNA-binding domain are cotransformed into a yeast reporter strain, and the resulting transformants are screened for those that express the reporter gene. For example, and not by way of limitation, the bait gene can be cloned into a vector such that it is translationally fused to the DNA encoding the DNA-binding domain of the GAL4 protein. These colonies are purified and the library plasmids responsible for reporter gene expression are isolated. DNA sequencing is then used to identify the proteins encoded by the library plasmids-

The present invention, thus generally described, will be understood more readily by reference to the following examples, which are provided by way of illustration and are not intended to be limiting of the present invention.

The examples below are provided to illustrate the subject invention. These examples are provided by way of illustration and are not included for the purpose of limiting the invention.

EXAMPLES

EXAMPLE I: cDNA Library Construction

10

15

20

25

30

35

cDNA library plates and clones originated from five cDNA libraries that were constructed by directional cloning. These are available through the Resource Center (http://www.rzpd.de) of the German Genome Project. In particular, the hfbr2 (human fetal brain; RZPD number DKFZp564) and hfkd2 (human fetal kidney; DKFZp566) libraries were generated using the Smart kit (Clontech), except that PCR was carried out with primers that contained uracil residues to permit directional cloning without restriction digestion and ligation, and were complementary with the pAMP1 (LifeTechnologies) cloning sites for directional cloning. The htes3 (human testes; DKFZp434), hutel (human uterus; DKFZp586) and hmcfl (human mammary carcinoma: DKFZp727) libraries are conventional (Gubler, U., Hoffman, B.J., (1983), A simple and very efficient method for generating cDNA libraries. Gene 25, 263-269), size-selected cDNA libraries. They are cloned into pSPORT1 (LifeTechnologies) via a NotI site which is introduced during reverse transcription downstream of the oligo dT primer and a SalI site that is introduced by the ligation of a adapters. The human mammary carcinoma library was constructed from MCF7 cells.

In a similar fashion, the hamy2 (human amygdala nucleus (inside the brain); RZPD number DKFZp7b1) and hmel2 (human melanoma; RZPD number DKFZp7b2) libraries have been generated using conventional approaches, emplying a NotI -dT V primer for first strand synthesis (GAGCGGCCGC(T)19V). After second strand synthesis, SalI adapters were ligated to the blunted cDNA. Then the cDNA was cut with NotI to generate SalI-NotI compatible ends at the 5' and 3' ends of the cDNA, respectively, to allow directional cloning. The cDNAs were then size selected on agarose gels in two dimensions and cloned into the pSPORT1 plasmid vector which had been pre-cut with SalI and NotI (LifeTechnologies). The

DNA was transformed into the DHLOB bacterial strain and single colonies were picked into 384well microtiter plates from the non-amplified library. The human melanoma library was constructed from MeWo cells, published by Kern, M.A., Helmbach, H., Artuc, 5 M., Karmann, D., Jurgovsky, K. and Schadendorf, D. (1997) Human melanoma cell lines selected in vitro displaying various levels of drug resistance against cisplatin, fotemustine, vindesine or etoposide: modulation of proto-oncogene expression. Anticancer Res. 17, 4359-4370-

The cDNA sequences of this application were first identified among the sequences comprising various libraries. Technology has advanced considerably since the first cDNA libraries were made. Many small variations in both chemicals and machinery have been instituted over time, and these have improved both the efficiency and safety of the process. Although the cDNAs could be obtained using an older procedure, the procedure presented in this application is exemplary of one currently being used by persons skilled in the art. For the purpose of providing an exemplary method, the mRNA isolation and cDNA library construction

20 described here is for the MCF-7 library (DKFZp727) from which the clones named DKFZphmcfl_xxyyxx were obtained.

The human cell line MCF-7 was grown in DMEM supplemented with 10% fetal calf serum until confluency. 3 X 108 cells were harvested with a cell scraper in PBS. Cells were lysed in buffer containing 0.5% NP-40 to leave the nuclei intact. The debris was pelleted by centrifugation at 15 000 x g for 10 minutes at 4 degrees Celsius. Proteins in the supernatant were degraded in presence of SDS and Proteinase K (30 minutes at 55 degrees Celsius). Precipitation of proteins was done in a Phenol/Chloroform extraction. RNA was precipitated from the aqueous phase with Na-acetate and Ethanol. Polyadenylated messages were isolated using @iagen Oligotex (@IAGEN. Hilden Germany).

First strand cDNA synthesis was accomplished using an oligo 35 (dT) primer which also contained an NotI restriction site. Second strand synthesis was performed using a combination of DNA polymerase I. E. coli ligase and RNase H. followed by the

addition of a SalI adaptor to the blunt ended cDNA. The SalI adapted, double-stranded cDNA was then digested with NotI restriction enzyme, and fractionated by size on an agarose gel-DNA of the appropriate size was cut from the gel and cast into a second gel in a 90° angle. After electrophoresis in the second dimension, cDNA of the appropriate size was cut from the gel. The agarose block was broken down with help of gelase. The cDNA was purified with help of two phenol extractions and an ethanol precipitation. The cDNA was ligated into Sall/NotI pre-digested pSportl vector (LifeTechnologies) and transformed into DHLOB bacteria.

The libraries were arrayed into 384-well microtiter plates and spotted on high density nylon membranes for hybridization analysis. All libraries have been arrayed into 384well microtiter plates and spotted on high density nylon membranes for hybridization analysis. The hamy2 Library consists of 121 384well plates comprising 46464 clones. The hmel2 library consists of 72 384well plates comprising 27648 clones. Filters and clones are available through the Resource Center of German Genome Project (http://www.RZPD.de). Whole library plates were distributed to the sequencing partners of the consortium for systematic ---sequencing ------

EXAMPLE II: Sequencing of cDNA Clones 25

10

15

20

30

All clones in the 384-well microtiter plates were sequenced from the 5' end. Sequencing was done preferentially using dye terminator chemistry (ABD or Amersham) on ABI automated DNA sequencers (ABI 377, Applied Biosystems), one partner used EMBL prototype instruments (Arakis) mainly with dye primer chemistry.

The resulting expressed sequence tag (EST) sequences ("rl . ESTs" = sequenced from 5'-end) were analysed for:

a) the lack of identical matches with known genes.

For this, the EST-sequence was blasted against the cDNA 35 consortiums own database and after that against public databases

and (with BLASTn and BLASTx against EMBL/EMBLNEW and assembled ESTs, please refer to EXAMPLE III: Bioinformatics analysis of full length cDNAs, for description and parameter settings). ESTs which were identical to known genes in more than 100 bp, with less than 2 mismatches, were excluded from further analysis.

b) the presence of an open reading frame

Open reading frames (ORFs) were detected with an tool developed by Munich Information Center for Protein Sequences (MIPS) called ORF-map. ORF-map visualises potential start and stop-codons. If an ORF without a stop codon was detected in a rl-EST, the sequence was processed further.

c) the presence of GC rich sequences

10

15

20

25

A script developed by MIPS computed the GC-content of the rl-sequence, which should be >40%. Writing similar scripts is within the ordinary skill of one in bioinformatics.

d) the lack of repeat structures

Repeats such as Alu, Line or CA-repeats were detected by blasting (BLASTn and BLASTx, please refer to EXAMPLE III: Bioinformatics analysis of full length cDNAs, for description and parameter settings) against a repeat-database compiled by MIPS. If a repeat was present within the rl-sequence, the sequence were not processed further.

Novel clones that met all criteria were identified to the sequencers, who then performed 3'-end sequencing of these clones. The resulting 3' ESTs ("sL ESTs" = sequenced from 3'-end) were checked for

a) the lack of matches with known genes in public databases, and sequences already generated by us.

This was done by blasting against EMBL/EMBLNEW and assembled EST (BLASTn and BLASTx, please refer to EXAMPLE III:

Bioinformatics analysis of full length cDNAs, for description and parameter settings).

b) the presence of polyadenylation signals.

5

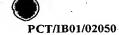
10

15

Again only clones matching the selection criteria were chosen to be sequenced completely by the sequencers. Clones were selected after the following criteria:

A very good ORF had at least one BLASTx match to other proteins. A "good ORF" should extend to the 3' end and be longer than ~40 codons. If the ORF started in the rl sequence, in front of the potential start codon, there should not exist too many competing start codons in frame with the ORF start codon and the start should match the Kozak consensus ATG. If the EST sequence was to short to decide according to the potential ORF, and there were only a few or no start codons in the sequence the GC content of the Sequence should be greater than 40%. The rl sequences needed not contain an polyA-tail at the 3' end. In addition, the results of the blasting against the assembled human ESTs could help in questionable cases to decide whether to stop or to continue. A hit against these ESTs was an indication to go further.

Clones passing the above-described screening were sequenced 20 in full. Sequencing was done preferentially using dye terminator chemistry (ABD or Amersham) on ABI automated DNA sequencers (ABI 377. Applied Biosystems), one partner used EMBL prototype instruments (Arakis) mainly with dye primer chemistry. Primer walking (Strauss et al., 1986, Specific-primer-directed DNA sequencing. Anal Biochem. 154, 353-360) was the preferred 25 sequencing strategy because of the lower redundancy possible compared to random shotgun (Messing, J., Crea, R., Seeburg, H.P. (1981) A system for shotgun DNA sequencing. Nucleic Acids Res. 9, 32-39) methods. Walking primers were generally designed using 30 software (e.g. Haas, S., Vingron, M., Poustka, A., Wiemann, S. (1998) Primer design in large-scale sequencing. Nucleic Acids Res. 26, 3006-3012; Schwager; C., Wiemann, S., Ansorge; W. (1995) GeneSkipper: integrated software environment for DNA sequence assembly and alignment. HUGO Genome Digest 2, 8-9) that permitted 35 complete automation of this usually time consuming process and helped in the parallel processing of large numbers of clones.



WO 01/98454 PCT/IB01/020
EXAMPLE III: Bioinformatics analysis of full length cDNAs

Each sequence obtained was compared on nucleotide level in a stepwise manner to sequences in EMBL/EMBLNEW, EMBL-EST, EMBL-STS using the BLASTn algorithm. Basic Local Alignment Search Tool (BLAST, Altschul S. F. (1993) J Mol Evol 35:290-300; Altschul, S. F. et al (1990) J Mol Biol 215:403-10) is used to search for local sequence alignments. BLAST produces alignments of both nucleotide (BLASTn) and amino acid sequences (BLASTp or BLASTx) to determine sequence similarity. BLAST is especially useful in determining exact matches or in identifying homologs, because of the local nature of the alignments. While it is useful for matches which do not contain gaps, it is inappropriate for performing motif-style searching. The fundamental unit of BLAST algorithm output is the High-scoring Segment Pair (HSP).

An HSP consists of two sequence fragments of arbitrary but equal lengths whose alignment is locally maximal and for which. the alignment BLAST approach is to look threshold or cut off score set by the user. BLAST looks for HSPs between a query sequence and a database sequence, to evaluate the statistical significance of any matches found, and to report only those which satisfy the user-selected threshold of significance. The parameter E establishes the statistically significant threshold for reporting database sequence matches. E is interpreted as the upper bound of the expected frequency of chance occurrence of an HSP (or set of HSPs) within the context of the entire database search. Any database sequence whose match satisfies E is reported in the program output. Parameter settings for the BLAST-operations (BLASTN 2-Dal9MP-WashU) described were: EMBL-EMBLNEW: H=O V=5 B=5 -filter seg; EMBL-EST: H=O E=le-lo B=500 V=500 -filter seg: EMBL-STS: H=0 V=5 B=5.

15

20

25

30

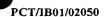
Search against EMBL/EMBLNEW was done to determine whether the cDNAs are already known, and also to find out whether the cDNAs are encoded by genomic sequences already sequenced and published/submitted to these databases.

35 Search against EMBL-EST was performed to get a first impression how abundant a particular cDNA would be and to get

WO 01/98454

10

. 15



information on tissue specificity (so-called "electronic Northern-Blot", e.g. some of the cDNAs derived of the testis library show only hits to ESTs also derived of testis libraries).

The cDNA-sequences were blasted against EMBL-STS to determine STS-sequence-match to the cDNA, thus providing a mapping information to the new cDNA.

The potential protein-sequences were generated automatically by a script searching for the longest open reading frame (ORF) in each of the three forward frames with a minimum length of 90 codons. Next, the automatically generated ORFs were translated into protein sequences. These protein sequences were searched the non redundant protein data set against PIR/SwissProt/Trembel/Tremblnew (BLASTP 2.Dal9MP-WashU, parameter setting: V=7 B=7 H=0 -filter seg). If the script generated more than one ORF, one ORF was chosen manually by the annotater according to the degree of similarity to known proteins, the location of the ORF in the cDNA, the length, the amino acid composition and the content of Prosite-Motifs.

Additionally there was a BLASTx (BLASTX 2.Dal9MP-WashU 20 redundant protein database comprising PIR/SWISSPROT/TREMBL/TREMBLNEW; parameter-settings ----matrix/home/data/blast/matrix/aa/BL0SUM62--H=0--V=5--B=5 -filter seg) search to find potential frame shift in the complementary cds of the cDNAs and to identify unspliced or partly spliced 25 cDNAs. The protein sequence was then transferred to the PEDANT system, in order to generate additional information on the new proteins. PEDANT (Protein Extraction, Description, and ANalysis Tool, Frishman, D. & Mewes, H.-W. (1997) PEDANTic, genome analysis. Trends in Genetics , 13, 415-416) is a platform developed at the Munich Information Center for Protein Sequences 30 (MIPS: Munich: Germany): which incorporates practically bioinformatics methods important for the functional structural characterisation of protein sequences. Computational methods used by PEDANT are:

WO 01/98454

FASTA

Very sensitive protein sequence database searches with estimates of statistical significance. Pearson W-R- (1990) Rapid and sensitive sequence comparison with FASTP and FASTA. Methods 5 Enzymol. 183, 63-98.

PCT/IB01/02050

BLAST2

Very sensitive protein sequence database searches with estimates of statistical significance. Altschul S.F., Gish W., Miller W., Myers E.W., and Lipman D.J. Basic local alignment search tool. Journal of Molecular Biology 215, 403-10.

PREDATOR

10

15

High-accuracy secondary structure prediction from single and multiple sequences. Frishman, D. and Argos, P. (1997) 75% accuracy in protein secondary structure prediction. Proteins, 27, 329-335. Frishman, D. and Argos, P. (1996) Incorporation of long-distance interactions in a secondary structure prediction algorithm. Prot. Eng. 9, 133-142.

STRIDE

Secondary structure assignment from atomic coordinates.

20 Frishman, D. and Argos, P.(1995) Knowledge-based secondary structure assignment. Proteins 23, 566-579.

CLUSTALW

Multiple sequence alignment. Thompson, J.D., Higgins, D.G. and Gibson, T.J. (1994) CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, positions-specific gap penalties and weight matrix choice. Nucleic Acids Research, 22:4673-4680.

TMAP

Transmembrane region prediction from multiply aligned

30 sequences. Persson, B. and Argos, P. (1994) Prediction of transmembrane segments in proteins utilising multiple sequence alignments. J. Mol. Biol. 237, 182-192.

WO 01/98454 ALOM2

Transmembrane region prediction from single sequences.

Klein, P., Kanehisa, M., and DeLisi, C. Prediction of protein function from sequence properties: A discriminant analysis of a database. Biochim. Biophys. Acta 787, 221-226 (1984). Version 2 by Dr. K. Nakai.

SIGNALP .

Signal peptide prediction Nielsen, H., Engelbrecht, J., Brunak, S., and von Heijne, G (1997). Identification of prokaryotic and eukaryotic signal peptides and prediction of their cleavage sites. Protein Engineering 10, 1-6.

SEG

Detection of low complexity regions in protein sequences.

Wootton, J.C., Federhen, S. (1993) Statistics of local complexity

in amino acid sequences and sequence databases. Computers &

Chemistry 17, 149-163.

COILS

Detection of coiled coils. Lupas, A., M. Van Dyke, and J.

Stock, "Predicting Coiled Coils from Protein Sequences." Science

(1991) 252, 1162-1164.

PROSEARCH

Detection of PROSITE protein sequence patterns. Kolakowski L.F. Jr., Leunissen.J.A.M., Smith.J.E. (1992) ProSearch: fast searching of protein sequences with regular expression patterns related to protein structure and function. Biotechniques 13, 919-921.

BLIMPS

Similarity searches against a database of ungapped blocks.

J.C. Wallace and Henikoff S., (1992) PATMAT: a searching and

extraction program for sequence, pattern and block queries and databases, CABIOS &, 249-254. Written by Bill Alford.

WO 01/98454 HMMER

Hidden Markov model software . Sonnhammer E.L.L., Eddy S.R., Durbin R. (1997) Pfam: A Comprehensive Database of Protein Families Based on Seed Alignments. Proteins 28, 405-420.

5 pI

10

15

Perl script that returns the amino acid composition, molecular weight, theoretical pI, and expected extinction coefficient of an amino acid sequence. By Fred Lindberg. The parameter-settings were as follows: known3d: score > 100; BLAST: E-value < 10; SCOP: <= 50 Alignments, E-Value < 0.0001; signalp: Y=0.7; untersucht vom N-Terminus her: 50 aa; funcat: E-value < 0.001; BLOCKS: <= 10 hits; BLIMPS: threshold 1100.0; COILS: threshold 0.95; SEG: threshold 20.0; BLAST in report: E-value < 0.001; PIR-KW, superfamilies, EC-Nummern in report: E-value < 0.00001; known3d in report: score > 120

The results of PEDANT analysis together with the results of the similarity searches constitute the basis for the structural and functional annotation of the cDNAs and the encoded proteins, as specified herein.

We claim:

An assemblage, comprising at least one nucleic acid molecule having the sequence of a clone selected from the group consisting of: amy2_l2g7; amy2_l2il; amy2_l3gl9; amy2_l6el4; amy2_24k15; amy2_2al3; amy2_2il7; fbr2_78dl8; fbr2_78el8; amy2_l2lm2; amy2_24b4; amy2_l2lfl9; tes3_l6b5; amy2_li24; amy2_ljl9; amy2_2bl9; amy2_7j5; amy2_l4b5; amy2_2ol3; fkd2_3kl; mel2_7gl4; mel2_l2jl ; mel2_7kl9; amy2_2c22; fbr2_78i2l; amy2_lln4; amy2_lcl2; amy2_lil; amy2_2f22; amy2_2gl2; fbr2_78cl2; tes3_10i16; tes3_3la10; amy2_10h17; amy2_10p7; amy2_12d7; 10 amy2_2fl&: tes3_11c22: tes3_11d21: tes3_29f24: tes3_31j20: tes3_5k22; Tes3_10n10; Tes3_1le17; Tes3_12d18 ; Tes3_141?; Tes3_15n14; Tes3_16p3; Tes3_19p12; Tes3_21k14; Tes3_22ill; Tes3_22124; tes3_2bq3; tes3_30pb; amy2_1ld2; amy2_12lo17; 15 amy2_lil4; amy2_24c8; fbr2_78d4; tes3_llal7; tes3_l7i2l; tes3_20hl2; tes3_7nl2; tes3_9el6; amy2_14ml6; tes3_18nl4; their complements; and variants thereof.

- 2. An assemblage, comprising at least one nucleic acid
 20 molecule having the sequence of a clone selected from the group consisting of: amy2_l2g7; amy2_l2il; amy2_l3gl9; amy2_lbel4; amy2_24kl5; amy2_2al3; amy2_2il7; amy2_l2lm2; amy2_24b4; amy2_l2lfl9; amy2_li24; amy2_ljl9; amy2_2bl9; amy2_7j5; amy2_l4b5; amy2_2ol3; amy2_2c22; amy2_lln4; amy2_lcl2; amy2_lil; amy2_lcl2; amy2_lil; amy2_lcl2; amy2_lil; amy2_lcl2; amy2
- 3. An assemblage, comprising at least one nucleic acid molecule having the sequence of a clone selected from the group 30 consisting of: fbr2_78dl8; fbr2_78el8; fbr2_78i2l; fbr2_78cl2; fbr2_78d4; their complements; and variants thereof.
 - 4. An assemblage, comprising at least one nucleic acid molecule having the sequence of a clone selected from the group consisting of: amy2_121m2; amy2_24b4; their complements; and variants thereof.

5. An assemblage, comprising at least one nucleic acid molecule having the sequence of a clone selected from the group consisting of: amy2_l2lfl9; tes3_lbb5; their complements; and variants thereof.

- 5 An assemblage comprising at least one nucleic acid molecule having the sequence of a clone selected from the group consisting of: amy2_li24; amy2_lj19; amy2_2b19; amy2_7j5; their complements; and variants thereof.
- 7. An assemblage, comprising at least one nucleic acid 10 molecule having the sequence of a clone selected from the group consisting of: amy2_14b5; amy2_2ol3; fkd2_3kl; mel2_7gl4; their complements; and variants thereof.
- 8. An assemblage, comprising at least one nucleic acid molecule having the sequence of a clone selected from the group consisting of mel2_7gl4; mel2_l2jl; mel2_7kl9; their complements; and variants thereof.
 - 9. An assemblage, comprising at least one nucleic acid molecule having the sequence of a clone selected from the group consisting of: amy2_2c22; fbr2_78i21; their complements; and variants thereof.

20

- 10. An assemblage, comprising at least one nucleic acid molecule having the sequence of a clone selected from the group consisting of: amy2_lln4; amy2_lil; amy2_2gl2; fbr2_7&cl2; tes3_l0ilb; tes3_3lal0; their complements; and variants thereof.
- 25 LL. An assemblage comprising at least one nucleic acid molecule having the sequence of a clone selected from the group consisting of: amy2_10h17; amy2_10p7; amy2_12d7; amy2_2f18; tes3_11c22; tes3_11d21; tes3_29f24; tes3_31j20; tes3_5k22; their complements; and variants thereof.
- 30 l2. An assemblage, comprising at least one nucleic acid molecule having the sequence of a clone selected from the group consisting of: tes3_16b5; tes3_10i16; tes3_31a10; tes3_11c22; tes3_11d21; tes3_29f24; tes3_31j20; tes3_5k22; Tes3_10n10; Tes3_11e17; Tes3_12d18; Tes3_1417; Tes3_15n14; Tes3_16p3;

Tes3_19pl2; Tes3_21kl4; Tes3_22ill; Tes3_22124; tes3_2bg3; tes3_30pb; tes3_1lal7; tes3_17i2l; tes3_20hl2; tes3_7nl2; tes3_9elb; their complements; and variants thereof.

- 13. An assemblage, comprising at least one nucleic acid molecule having the sequence of a clone selected from the group consisting of: amy2_lld2; amy2_l2lol7; amy2_lil4; amy2_24c8; fbr2_78d4; tes3_llal7; tes3_l7i2l; tes3_20hl2; tes3_7nl2; tes3_9el6; their complements; and variants thereof.
- 10 l4. An assemblage, comprising at least one nucleic acid molecule having the sequence of a clone selected from the group consisting of: amy2_14mlb; tes3_18nl4; amy2_1cl2; amy2_2f22; their complements; and variants thereof.
- 15. A nucleic acid molecule comprising a nucleotide 15 sequence of the clone fkd2_3k1.
- A computer readable medium, comprising in electronic form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: amy2_12g7; amy2_12il; amy2_13g19; amy2_16e14; amy2_24k15; amy2_2a13; amy2_2i17; fbr2_78dl8; fbr2_78el8; amy2_121m2; amy2_24b4; amy2_121f19; 20 tes3_16b5; amy2_1i24; amy2_1j19; amy2_2b19; amy2_7j5; amy2_14b5; amy2_2ol3: fkd2_3kl: mel2_7gl4: mel2_l2jl : mel2_7kl9: amy2_2c22: fbr2_78i2l; amy2_lln4; amy2_lcl2; amy2_lil; amy2_2f22; amy2_2gl2; fbr2_78cl2; tes3_l0ilb; tes3_3lalO; amy2_l0hl7; amy2_l0p7; 25 amy2_12d7; amy2_2f18; tes3_11c22; tes3_11d21; tes3_29f24; tes3_31j20; tes3_5k22; Tes3_10n10; Tes3_11e17; Tes3_12d18; Tes3_1417; Tes3_15n14; Tes3_16p3; Tes3_19p12; Tes3_21k14; Tes3_22ill; Tes3_22124; tes3_26g3; tes3_30p6; amy2_lld2; amy2_l2lol7; amy2_lil4; amy2_24c8; fbr2_78d4; tes3_llal7; tes3_17i21; tes3_20h12; tes3_7n12; tes3_9e16; amy2_14m16; 30 tes3_lanl4; their complements; and variants thereof.
 - 17. A computer readable medium, comprising in electronic form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: amy2_12g7; amy2_12il; amy2_13g19; amy2_16e14; amy2_24k15; amy2_2al3; amy2_2il7;

amy2_l2lm2; amy2_24b4; amy2_l2lfl9; amy2_li24; amy2_ljl9; amy2_2bl9; amy2_7j5; amy2_l4b5; amy2_2ol3; amy2_2c22; amy2_lln4; amy2_lcl2; amy2_lil; amy2_2f22; amy2_2gl2; amy2_l0hl7; amy2_l0p7; amy2_l2d7; amy2_2fl8; amy2_lld2; amy2_l2lol7; amy2_lil4; amy2_24c8; their complements; and variants thereof.

18. A computer readable medium, comprising in electronic form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: fbr2_78dl8; fbr2_78el8; fbr2_78i2l; fbr2_78cl2; fbr2_78d4; their complements; and variants thereof.

10

25

30

- 19. A computer readable medium, comprising in electronic form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: amy2_121m2; amy2_24b4; their complements; and variants thereof.
- 15 20. A computer readable medium, comprising in electronic form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: amy2_121f19; tes3_1665; their complements; and variants thereof.
- 21. A computer readable medium, comprising in electronic 20 form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: amy2_li24; amy2_lj19; amy2_2b19; amy2_7j5; their complements; and variants thereof.
 - 22. A computer readable medium, comprising in electronic form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: amy2_14b5; amy2_2ol3; fkd2_3kl; mel2_7gl4; their complements; and variants thereof.
 - 23. A computer readable medium, comprising in electronic form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: mel2_l2jl; mel2_7kl9; their complements; and variants thereof.
 - 24. A computer readable medium, comprising in electronic form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: amy2_2c22; fbr2_7&i21; their complements; and variants thereof.

25. A computer readable medium, comprising in electronic form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: amy2_lln4; amy2_lil; amy2_2gl2; fbr2_7&cl2; tes3_l0ilb; tes3_3lal0; their complements; and variants thereof.

26. A computer readable medium, comprising in electronic form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: amy2_l0hl7; amy2_l0p7; amy2_l2d7; amy2_2fl8; tes3_llc22; tes3_lld2l; tes3_29f24; tes3_3lj20; tes3_5k22; their complements; and variants thereof.

10

15

20

25

30

- 27. A computer readable medium; comprising in electronic form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: tes3_lbb5; tes3_l0ilb; tes3_3lal0; tes3_llc22; tes3_lld2l; tes3_29f24; tes3_3lj20; tes3_5k22; Tes3_l0nl0; Tes3_llel7; Tes3_l2dl8; Tes3_l417; Tes3_l5nl4; Tes3_lbp3; Tes3_l9pl2; Tes3_2lkl4; Tes3_22ill; Tes3_22l24; tes3_2bg3; tes3_30pb; tes3_llal7; tes3_l7i2l; tes3_20hl2; tes3_7nl2; tes3_9elb; their complements; and variants thereof.
- 28. A computer readable medium; comprising in electronic form at least one nucleic acid or protein sequence of a clone -selected from the group consisting of: amy2_lld2; amy2_l2lol7; amy2_lil4; amy2_24c8; fbr2_78d4; tes3_llal7; tes3_l7i2l; tes3_20hl2; tes3_7nl2; tes3_9el6; their complements; and variants thereof.
 - 29. A computer readable medium, comprising in electronic form at least one nucleic acid or protein sequence of a clone selected from the group consisting of: amy2_14mlb; tes3_18nl4; amy2_1cl2; amy2_2f22; their complements; and variants thereof.
 - 30. A computer readable medium, comprising in electronic form a nucleic acid or protein sequence of the clone fkd2_3k1.